



**HAL**  
open science

## Distance Learning for Analog Methods

Paul Platzer, Arthur Avenas, Bertrand Chapron, Lucas Drumetz, Alexis Mouche, Pierre Tandeo, Léo Vinour

► **To cite this version:**

Paul Platzer, Arthur Avenas, Bertrand Chapron, Lucas Drumetz, Alexis Mouche, et al.. Distance Learning for Analog Methods. 2024. hal-04841334

**HAL Id: hal-04841334**

**<https://hal.science/hal-04841334v1>**

Preprint submitted on 16 Dec 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

## Distance Learning for Analog Methods

Paul Platzer<sup>a b</sup>, Arthur Avenas<sup>a b c</sup>, Bertrand Chapron<sup>a b</sup>, Lucas Drumetz<sup>b d</sup>,  
Alexis Mouche<sup>a</sup>, Pierre Tandeo<sup>b d e</sup>, Léo Vinour<sup>a f</sup>

<sup>a</sup> *Laboratoire d'Océanographie Physique et Spatiale, Univ. Brest/Ifremer/CNRS/IRD, F-29280  
Plouzané, France*

<sup>b</sup> *Odyssey, Inria/IMT/CNRS, F-29280 Plouzané, France*

<sup>c</sup> *European Space Agency, Esrin, Italy*

<sup>d</sup> *IMT Atlantique, Lab-STICC, UMR CNRS 6285, 29238, Brest, France*

<sup>e</sup> *RIKEN Cluster for Pioneering Research, Kobe, 650-0047, Japan*

<sup>f</sup> *France Énergies Marines, Plouzané, France*

12 **ABSTRACT:** Analogs are similar states of a system, occurring at remote times within independent  
13 numerical simulations or previous observations. This concept has been developed in atmospheric  
14 sciences, and was further used in atmospheric and ocean sciences for forecasting, downscaling,  
15 upscaling, extreme event attribution, and many other applications. The distance used to find and  
16 rate analogs is a key feature of analog methods. Most studies are based on the Euclidean distance or  
17 other pre-defined metrics. In this investigation, we leverage distance learning algorithms originally  
18 designed for classification and regression and adapt them for statistical forecasting objectives, using  
19 in particular the continuous-ranked probability score as a loss function. Our algorithm allows to  
20 jointly optimize three key hyperparameters of analog methods: the feature space, the distance,  
21 and the number of analogs used. In particular, this algorithm allows to reduce the feature space  
22 dimension while keeping analog ensemble performances as high as possible, a key requirement  
23 for small and medium-sized datasets. We test our algorithm on an idealized chaotic system and on  
24 a small-size tropical cyclone dataset from meteorological agencies. Our experiments suggest that  
25 the optimal distance strongly depends on the forecast horizon and the number of available data,  
26 and that our algorithm allows for reasonable performances of analog ensemble methods even for  
27 small-size datasets. Our approach is not limited to forecasting and can assist the search for optimal  
28 hyperparameters of any analog method, enhancing exploration possibilities and improving overall  
29 performances.

30 SIGNIFICANCE STATEMENT: “History repeats itself.” Today’s weather is likely to be remi-  
31 niscent of past, already observed weather. This is the notion of “analog” weather. Analogs were  
32 introduced to study the atmosphere, but were also used recently to study the ocean. They are used  
33 in many applications, including forecasts, or to estimate whether extreme events are influenced by  
34 climate change or not. To decide whether two distant-time images of the atmosphere (or ocean)  
35 are “analogs” of each other, one must define a similarity criterion, the “distance”. The definition  
36 of the distance depends on the chosen application. This research aims at providing an algorithm to  
37 systematically tune the distance used in the definition of analogs, depending on the application.

38 *This Work has been submitted to Monthly Weather Review. Copyright in this Work may be*  
39 *transferred without further notice.*

## 40 **1. Introduction**

41 Analog methods rely on the search for neighbours of any given query in a database, also called  
42 ”catalog” or ”library”, such as a reanalysis (*e.g.* ERA5, Hersbach et al. 2020) or an ensemble  
43 simulation (*e.g.* CMIP5, Taylor et al. 2012). They have been used in a wide range of applications  
44 in atmosphere and ocean sciences, including downscaling (Zorita and Von Storch 1999), upscaling  
45 (Yiou et al. 2014), ensemble forecasts (Delle Monache et al. 2013; Yiou 2014), tropical cyclone  
46 forecasting (Neumann and Hope 1972), extreme events detection and attribution (Jézéquel et al.  
47 2018), importance sampling (Yiou and Jézéquel 2020), data assimilation (Tandeo et al. 2015;  
48 Lguensat et al. 2017), and interpolation (Zhen et al. 2020).

49 In atmospheric science, the idea of looking in observational archives to find similar states dates  
50 back to the concept of “points of symmetry” used in weather forecasting at least since Weickmann  
51 (1924) and then by Krick (1942) and Elliott (1943). The term “analogs” can be attributed to the  
52 study of Lorenz (1969) on atmospheric predictability. This term refers to methods and concepts  
53 developed in atmospheric science and now also used in ocean science (*e.g.* Le Bras et al. 2024),  
54 while other terms such as “neighbours” are more general and may refer to other concepts from  
55 dynamical systems (Lucarini et al. 2016) or machine learning (Peterson 2009). In this paper,  
56 “analog methods” explicitly refers to the use of neighbours-based methods in the specific context  
57 of atmospheric and ocean science.

58 The popularity of analog methods is probably due to their simplicity of interpretation and  
59 implementation. One simply needs to define a feature space, a distance (and kernel if weighting  
60 is applied), and to choose a number of analogs to be used. Then, one can perform any statistical  
61 task (*e.g.* averaging, linear regression, ensemble forecast) on the analog ensemble. These 3 key  
62 hyperparameters of analog methods (feature space, distance, and number of analogs used) are often  
63 directly chosen by the user depending on the application. Therefore, experience with using analog  
64 methods is required to make suitable hyperparameters choices, which prohibits an even wider use  
65 of analog methods in the atmospheric and ocean science community. Moreover, the sensitivity of  
66 analog methods to such choices is usually little or not explored unless it is the sole purpose of the  
67 study.

68 Empirical rules for the choice of these hyperparameters are as follows. First, the choice of features  
69 is based on the user’s understanding of the physical system and on the objective task. Knowledge of  
70 dimensionality issues for analog methods in the case of limited-size datasets (Van den Dool 1994;  
71 Nicolis 1998; Platzer et al. 2021a) encourages one to reduce the number of features used to define  
72 analogs. Empirical orthogonal functions (EOFs, see *e.g.* Lorenz 1956) are commonly chosen as  
73 a dimension reduction technique to be applied before using analog methods (*e.g.* Benestad 2010).  
74 However EOFs were not developed for the purpose of analog methods, and are therefore likely to  
75 be sub-optimal.

76 Second, once features have been defined, one has to choose a distance between many available  
77 options (see, *e.g.*, the ones explored by Toth 1991; Matulla et al. 2008). Note that even after  
78 choosing a distance family (*e.g.* Euclidian  $l_2$  vs. Manhattan  $l_1$ ), an infinite number of variations of  
79 distance definitions are possible within each family, for instance by giving different weights to each  
80 feature, or by performing any one-to-one transformation of the feature vectors before computing  
81 the distance. This allows for a lot of explorations in the definition of the distance, so much that  
82 one usually cannot afford to perform by hand. As a consequence, the Euclidean distance in its  
83 simplest form is usually chosen for simplicity of implementation. Note also that the choice of  
84 distance only matters in the case of finite-sized datasets, as all distance functions<sup>1</sup> are equivalent  
85 in finite dimension (see, *e.g.*, appendix A in Platzer et al. 2021a). Therefore, the choice of distance  
86 is expected to be more important when using datasets of small or medium size.

---

<sup>1</sup>A distance function “ $\text{dist}(\cdot, \cdot)$ ” should be positive ( $\text{dist}[x, x'] > 0$ ), definite ( $\text{dist}[x, x'] = 0$  iff.  $x = x'$ ), symmetric ( $\text{dist}[x, x'] = \text{dist}[x', x]$ ), and satisfy the triangle inequality ( $\text{dist}[x, x''] \leq \text{dist}[x, x'] + \text{dist}[x', x'']$ ).

87 Third and finally, a near-optimal value of number of analogs is generally searched by trial-  
88 and-error, typically between  $\mathcal{O}(10) - \mathcal{O}(200)$ . If one requires a probabilistic estimation with an  
89 ensemble of analogs, one will need a large number of analogs, while for a deterministic estimation  
90 one analog can be enough in theory. However, even for deterministic estimates, the bias-variance  
91 trade-off rule (Hastie 2009) encourages one to use several analogs. A rule-of-thumb can be used  
92 by setting an upper-bound for the typical analog-to-target distance. Indeed, Platzer et al. (2021a)  
93 showed that the distance  $r_k$  to the  $k$ -th analog (ranked by growing distance) scales as  $r_k \sim (k/N)^{1/d}$   
94 where  $N$  is the number of independent samples in the catalog, and  $d$  is the local dimension, the latter  
95 being upper-bounded by  $p$ , the number of features in the distance definition. Setting the condition  
96  $r_k < 0.2$  (the value 0.2 is arbitrary), we have a sufficient rule for  $k$  which is  $k < N \times 0.2^p$ . Using  
97 such a rule for the maximum number of analogs to be used allows to have analog-to-target distances  
98 that do not exceed roughly 20% of the typical distance between two points chosen randomly in  
99 the dataset. However, this is only an approximate statistical upper-bound, and finding the optimal  
100 number of analogs for a given objective task usually requires a lot of testing efforts.

101 A typical meteorological system that was extensively studied with analog methods is the tropical  
102 cyclone (TC). In particular, analog methods have been used for the forecast of tropical cyclone  
103 tracks (*e.g.* Fraedrich et al. 2003; Langmack et al. 2012; Bonnardot et al. 2019; Bessafi et al.  
104 2002) and intensities (*e.g.* Fetanat et al. 2013; Elsberry and Tsai 2014; Tsai and Elsberry 2014;  
105 Chen et al. 2016; Alessandrini et al. 2018; Tsai and Elsberry 2019; Bonnardot et al. 2019; Lewis  
106 et al. 2021). The tropical cyclone dynamical parameters are compiled in global databases called  
107 *best-tracks* (IBTrACS, Knapp et al. 2010). However, because of their underlying extreme ocean  
108 and atmospheric states, the crucial physical parameters of these phenomenon are still today largely  
109 undersampled in space and time. Recent studies also suggested that high-resolution observational  
110 data, *e.g.* from satellite synthetic aperture radar, is required when studying this dynamical system  
111 (Avenas et al. 2023, 2024b). Thanks to new satellite missions and new observation strategies  
112 along with increased performances of the sensors, both the spatio-temporal resolution (Jackson  
113 et al. 2021) and the quality (Combot et al. 2020) of these data to capture the storm wind structure  
114 are greatly increasing during these years, soon providing a dataset of a reasonable size, although  
115 still limited, to be combined with statistical approaches such as analogs. However, such methods

116 will face the issues of small-size dataset and dimensionality mentioned earlier, and will therefore  
117 require a fine-tuning of the number of features and of the distance used.

118 Several authors have developed methodologies for the optimization of analog methods' hyper-  
119 parameters. The kernels of Zhao and Giannakis (2016) are adapted to the local dynamics of  
120 the system under study and therefore provide optimal forecasts in the limit of large catalog size,  
121 as later demonstrated by Alexander and Giannakis (2020) using reproducing kernel Hilbert space  
122 theory. However, these methods do not tackle explicitly dimension reduction and small-size dataset  
123 issues. Also, they are targeted at forecasting of dynamical systems, while many other tasks can  
124 be performed with analog ensemble, and also require the optimisation of hyperparameters. Mc-  
125 Dermott and Wikle (2016) have given a Bayesian formulation of analog forecasting, allowing for  
126 optimisation of parameters through log-likelihood minimization. This approach yields a precise  
127 probabilistic formulation, and is flexible. Using a similar approach, Horton et al. (2017) perform  
128 optimization of analog methods through genetic algorithms, which allows to find global optimum  
129 and search through a wide range and number of hyperparameters. However, these two last methods  
130 require a thorough definition of hyperparameters to be optimised, while a more unifying framework  
131 would allow for a simpler representation and algorithmic implementation of the optimization of  
132 analog methods' hyperparameters.

133 Some authors have especially focused on the sensitivity of analog methods to the choice of  
134 distance. Toth (1991) compared nine different distances for the purpose of forecasting 700mb  
135 geopotential fields, and found some superiority for the root-mean-square difference between gra-  
136 dients of geopotential. Matulla et al. (2008) did a similar exercise for the purpose of precipitation  
137 downscaling, using five different distances and time-delayed embeddings in the space of EOFs.  
138 Authors jointly studied the effect of truncation in EOF space, the length of time-embedding, and the  
139 choice of distance. These studies allow for a detailed analysis of the performance of each distance  
140 choice for a specified task. However, as pointed by the authors of these two studies, they cannot  
141 generalize to other tasks than the one studied (forecasting for the first study and downscaling for the  
142 second). Moreover, one can generally not afford to conduct such tedious studies when designing a  
143 given analog method.

144 Other authors have introduced advanced methodologies for optimizing distances used in analog  
145 methods although this was not the direct topic of the study. This is the case of the tropical-cyclone

146 intensity analog forecasting scheme of Alessandrini et al. (2018). To find the best variables in more  
147 than 60 possible choices, the authors first searched for the best variable used alone, and then added  
148 new variables one-by-one, performing grid-search to find the best added variable and to define the  
149 respective weights of each new variable in the definition of the distance. Although it may prove  
150 efficient, the final solution given by such an algorithm is likely to be sub-optimal. Indeed, the fact  
151 that one predictor variable performs best alone does not guarantee that it should be retained when  
152 using several variables altogether. Furthermore, grid-search is a computationally intensive method  
153 to search for an optimum.

154 Fraedrich and Rückert (1998) optimized the distance used in analog methods by acting iteratively  
155 on weights given to coordinates used in Euclidean distance. However, by defining their own  
156 optimization step-rule, which was not studied elsewhere in the literature, the authors deprive  
157 themselves from the wealth of knowledge available for other well-known methods such as gradient  
158 descent. Indeed, the latter has well-established convergence properties, efficient algorithms (*e.g.*  
159 Kingma and Ba 2014), and has been used extensively in other fields of research. Furthermore,  
160 the algorithm of Fraedrich and Rückert (1998) is limited to optimizing weights, and a natural  
161 generalization would be the optimization of any linear transformation of the feature vector. Finally,  
162 the algorithm of Fraedrich and Rückert (1998) does not allow to perform dimension reduction,  
163 which is especially important when using analog methods with limited-size datasets. On the  
164 contrary, gradient-descent algorithms include regularization techniques which can help perform  
165 dimension reduction through feature selection.

166 In this paper, we focus on optimizing the distance used in analog methods. In particular, we  
167 leverage advances made in the field of machine learning through what is called “distance learning”  
168 or “metric learning” (Bellet et al. 2022). To our knowledge, distance learning algorithms have  
169 not been used before in analog methods for atmospheric and ocean sciences. We modified the  
170 ”Metric Learning for Kernel Regression” algorithm (Weinberger and Tesauro 2007) by fixing the  
171 number of neighbors (here called “analog”) and adding a regularization term, just as Yang et al.  
172 (2012) modified the “Neighborhood Component Analysis” algorithm (Goldberger et al. 2004).  
173 Also, we adapted these algorithms to allow for probabilistic estimates such as ensembles, widely  
174 used in the context of atmospheric and ocean sciences. To our knowledge, this is the first time  
175 that a metric learning algorithm is adapted to probabilistic forecasts. For this purpose, we replace



176 the mean squared error in the loss function by the continuous ranked probability score (CRPS,  
177 see Hersbach 2000), a score which is well suited for ensemble-based predictions such as the ones  
178 provided by analog methods. Most metric learning algorithms, including the one presented here,  
179 aim at optimizing a Mahalanobis-type distance, which can be understood as a Euclidean distance  
180 computed after normalization with a positive semi-definite matrix, or equivalently after applying  
181 a linear transformation to the data. Note that what is called here “Mahalanobis-type distance”  
182 is more general than the “Mahalanobis distance” *per se* in which the data is renormalized by its  
183 covariance matrix. Here, the normalization matrix (or, equivalently, the linear transformation) is  
184 not set in advance but must rather be optimized.

185 Note that optimizing the distance *includes* optimizing the feature space: our algorithm allows  
186 to find optimal linear combination of any set of features, imposing sparsity if needed to perform  
187 dimension reduction, possibly removing irrelevant features. Furthermore, by adjusting the scale  
188 of the distance, we are able to adjust the bandwidth of analog methods, and therefore the number  
189 of analogs used. By adapting a simple, well-known algorithm from distance learning, we are  
190 therefore able to jointly optimize three key hyperparameters of analog methods: the feature space,  
191 the distance, and the number of analogs used. By basing our algorithm on one simple, well-known  
192 gradient-descent rule, we allow it to be applied to any existing analog methodology, in a unifying  
193 framework that simplifies the search for hyperparameters of analog methods and enables better  
194 overall performances.

195 The three-variable dynamical system of Lorenz (1963) is widely used to mimic atmospheric  
196 and oceanic systems because it is chaotic. It was developed as a simplified, low-order model for  
197 atmospheric convection, by focusing on first-order large-scale Fourier modes. We will use it as a  
198 well-known toy model to test the effect of our distance optimization algorithm on analog forecasting,  
199 allowing to explore in particular the effect of forecast horizon and catalog size, as well as comparing  
200 the results of our algorithm when using two different loss functions: the mean-squared error or the  
201 continuous-ranked probability score (CRPS).

202 We also build a sub-dataset based on IBTrACS tropical cyclone data, used here to show that  
203 the proposed algorithm allows to select relevant variables for intensity forecasting. First, this  
204 should demonstrate that our algorithm can be used on small-size datasets, which is an interesting  
205 properties for tropical cyclone (TC) studies. Second, this should show that our algorithm allows

206 to perform dimension reduction by removing variables, which is a very important requirement for  
207 analog methods, as they are sensitive to the curse of dimensionality.

208 This paper is organized as follows. Basics of analog methods are recalled in section 2, before  
209 introducing the proposed algorithm and associated gradients. Experiments on the Lorenz system  
210 and on IBTrACS tropical cyclone data are reported in section 3 and 4. Conclusion and perspectives  
211 are drawn in section 5.

## 212 **2. Algorithm**

### 213 *a. Analog methods*

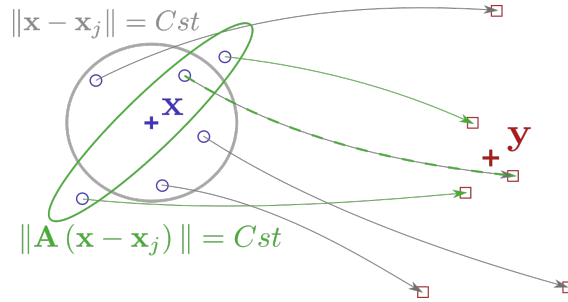
214 We assume that analogs of a query (or target)  $\mathbf{x}$  are sought for in a catalog  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  of  $N$   
215 independent  $d$ -dimensional vectors. The vectors  $\mathbf{x}_j$  can be, for instance, the values of any variable  
216 (temperature, geopotential height) on a lon-lat grid, at a given time. The idea of analog methods  
217 is to search for analogues of a “common” situation  $\mathbf{x}$ , in order to predict the associated output  
218 variable  $\mathbf{y}$ , such as the precipitation at a given station.

219 In general, analogs are given weights which are increasing functions of the similarity between  $\mathbf{x}$   
220 and  $\mathbf{x}_j$ . Therefore, a statistical estimate of  $\mathbf{y}$  is given by the weighted ensemble of values  $\{\mathbf{y}_j, j \in$   
221  $I(\mathbf{x})\}$  associated with the analogs. From there on, one can use either pure empirical ensemble  
222 prediction (Yiou 2014), or sample average and standard deviation assuming a Gaussian probability  
223 distribution for  $\mathbf{y}$ , or weighted linear regression computed on the analog sample  $\{(\mathbf{x}_j, \mathbf{y}_j), j \in I(\mathbf{x})\}$   
224 (Platzer et al. 2021b).

### 225 *b. Gradient of MSE of analog ensemble average*

226 Our algorithm finds an optimal linear transformation of explanatory variables  $\mathbf{x}$  before searching  
227 for analogs. The advantage of this algorithm is the simple computation of the gradient of the loss  
228 function. This is made possible by the use of Gaussian weights and a Mahalanobis-like distance.  
229 In our case, we will assume that the weights are of the following form:

$$p(\mathbf{x}_j | \mathbf{x}) \propto \exp\left(-\|\mathbf{A}(\mathbf{x} - \mathbf{x}_j)\|^2\right) \quad (1)$$



236 FIG. 1. Illustration of the process of making a linear transformation before applying analog ensemble methods.  
 237 On the left are shown the input variables (also called features), with target value  $\mathbf{x}$  (blue cross) surrounded by  
 238 potential analogs  $\mathbf{x}_j$  (blue circles). On the right are shown the output variables, with true output  $\mathbf{y}$  (red cross)  
 239 and outputs  $\mathbf{y}_j$  (red squares) associated with analogs  $\mathbf{x}_j$ . The selection of analogs is performed in input space  
 240 (or feature space), where the grey circle corresponds to a level-curve of constant Euclidean distance to the target  
 241 value  $\mathbf{x}$ , and the green circle corresponds to constant distance to  $\mathbf{x}$  after applying a linear transformation  $\mathbf{A}$ .

230 where  $\mathbf{A}$  is a  $p \times d$ -matrix with  $p \leq d$ , and  $\|\cdot\|$  is the Euclidean norm. After normalization over  
 231 the set of analogs  $I(\mathbf{x})$ , we find that each analog is given the empirical probability:

$$p(\mathbf{x}_j|\mathbf{x}) = \frac{\exp(-\|\mathbf{A}(\mathbf{x} - \mathbf{x}_j)\|^2)}{\sum_{k \in I(\mathbf{x})} \exp(-\|\mathbf{A}(\mathbf{x} - \mathbf{x}_k)\|^2)}. \quad (2)$$

232 The procedure of applying a linear transformation before selecting and weighing analogs is  
 233 illustrated in Fig. 1. In this example, applying the matrix  $\mathbf{A}$  allows to select analogs  $\mathbf{x}_j$  with  
 234 associated outputs  $\mathbf{y}_j$  that are closer to the true output  $\mathbf{y}$ , compared to the result of using the  
 235 original distances (with  $\mathbf{A}$  replaced by the identity matrix).

242 The objective is now to find the optimal  $p \times d$ -matrix  $\mathbf{A}$ , which is a linear transformation of  
 243 vectors  $\mathbf{x}$  to a space in which the analogs are sought for using the simple Euclidean distance. This  
 244 is equivalent to finding the optimal Mahalanobis-type distance  $\text{dist}(\mathbf{x}, \mathbf{x}_j) = (\mathbf{x} - \mathbf{x}_j)^T \mathbf{A}^T \mathbf{A} (\mathbf{x} - \mathbf{x}_j)$   
 245 (where  $T$ -subscript stands for “transpose”) with positive symmetric matrix  $\mathbf{A}^T \mathbf{A}$ .

246 It should be stressed again here that our use of the expression “Mahalanobis-type distance” differs  
 247 from another common use, where the symmetric, positive, semi-definite matrix  $\mathbf{A}^T \mathbf{A}$  is replaced  
 248 by the covariance of the data (McLachlan 1999), which would here be the catalog. Instead, we  
 249 optimize the matrix  $\mathbf{A}$ , and therefore we say that we optimize the distance within the family  
 250 of Mahalanobis-type distances. Note that some distance learning algorithms choose to directly

251 estimate the symmetric, positive, semi-definite matrix, which can sometimes allow to define convex  
 252 optimization algorithms (Globerson and Roweis 2005). However, optimizing  $\mathbf{A}$  has the advantage  
 253 to ease the interpretation, in particular in the context of dimension reduction.

254 To keep the analogy with the algorithm of Goldberger et al. (2004), we optimize the mean-  
 255 square error (MSE) of the analog prediction from the catalog, using a leave-one-out procedure.  
 256 We compute the square error of analog average prediction of  $\mathbf{y}_i$  from the truncated catalog  
 257  $\{\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_N\}$ . For simplicity, we note  $I(i) = I(\mathbf{x}_i)$  the indices of analogs of  $\mathbf{x}_i$  in  
 258 the truncated catalog, and  $p(j|i) = p(\mathbf{x}_j|\mathbf{x}_i)$  is the probability given to value  $\mathbf{y}_j$  associated with  
 259 analog  $\mathbf{x}_j$  of target  $\mathbf{x}_i$  using the truncated catalog. Note that  $p(j|i) \neq p(i|j)$  is not symmetric,  
 260 because of the normalization factor in Eq. (2). The MSE can then be written:

$$\text{MSE}(\mathbf{A}) = \frac{1}{N} \sum_{i=1}^N (\hat{\mathbf{y}}_i - \mathbf{y}_i)^T (\hat{\mathbf{y}}_i - \mathbf{y}_i), \quad (3)$$

261 where we use the notation  $\hat{\mathbf{y}}_i = \sum_{j \in I(i)} p(j|i) \mathbf{y}_j$  for the average analog prediction of  $\mathbf{y}_i$ . This MSE  
 262 has the following gradient:

$$\frac{\partial \text{MSE}}{\partial \mathbf{A}} = \frac{2}{N} \sum_{i=1}^N (\hat{\mathbf{y}}_i - \mathbf{y}_i)^T \sum_{k \in I(i)} \frac{\partial p(k|i)}{\partial \mathbf{A}} \mathbf{y}_k. \quad (4)$$

263 Note that what we call ‘‘gradient’’ for simplicity is actually a ‘‘sub-gradient’’, because we neglect  
 264 the discontinuities of  $\text{MSE}(\mathbf{A})$  due to changes in the the lists of analog indices  $I(i)$  with the change  
 265 of  $\mathbf{A}$ . The computation of the sub-gradient however is simpler than the true gradient, and sufficient  
 266 for convergence (Held et al. 1974). In the following we therefore only use the word ‘‘gradient’’.

267 Since we use exponential weights, the gradient can be computed easily. Indeed, the gradient of  
 268  $p(k|i)$  with respect to matrix  $\mathbf{A}$ , which is a matrix of the same shape as  $\mathbf{A}$ , is given by:

$$\frac{\partial p(k|i)}{\partial \mathbf{A}} = \mathbf{A} p(k|i) \left( \sum_{l \in I(i)} p(l|i) \mathbf{x}_{il} \mathbf{x}_{il}^T - \mathbf{x}_{ik} \mathbf{x}_{ik}^T \right), \quad (5)$$

269 where we use the notation  $\mathbf{x}_{ik} = \mathbf{x}_k - \mathbf{x}_i$ .

270 After simplifications, similar to the ones of Goldberger et al. (2004), the gradient of the MSE  
 271 can be rewritten:

$$\frac{\partial \text{MSE}}{\partial \mathbf{A}} = \mathbf{A} \frac{2}{N} \sum_{i=1}^N (\hat{\mathbf{y}}_i - \mathbf{y}_i)^T \sum_{j \in I(i)} p(j|i) (\hat{\mathbf{y}}_i - \mathbf{y}_j) \mathbf{x}_{ij} \mathbf{x}_{ij}^T. \quad (6)$$

272 where there are two independent products inside the sum: one *dot* product between vectors  $(\hat{\mathbf{y}}_i - \mathbf{y}_i)$   
 273 and  $(\hat{\mathbf{y}}_i - \mathbf{y}_j)$ , and a *tensor* product between the vector  $\mathbf{x}_{ij}$  and itself. Applying the matrix  $\mathbf{A}$  of  
 274 shape  $(p, d)$  to the tensor product of shape  $(d, d)$  gives a matrix of the same size as  $\mathbf{A}$ , which is the  
 275 size of the gradient.

276 The advantage of expressing the gradient in this way is that many terms are already computed  
 277 during the analog average prediction step. Therefore, to compute the gradient of the MSE on the  
 278 catalog, one must simply apply matrix  $\mathbf{A}$  to a product of partly pre-computed terms.

### 279 *c. Gradient of CRPS of analog ensemble*

280 In atmospheric and ocean science, estimation of uncertainty is a key feature of any estimation  
 281 algorithm. Also, the strength of analog methods is the cheap generation of ensembles. Therefore,  
 282 using the MSE to assess analog methods is too restrictive and we propose to use the continuous-  
 283 ranked probability score (CRPS), widely used in atmospheric sciences (Hersbach 2000).

284 The CRPS of a one-dimensional probabilistic forecast with cumulative probability distribution  
 285  $F$ , compared to a true scalar outcome  $y$ , can be expressed as :

$$\text{CRPS}(F, y) = \int_{-\infty}^{\infty} [F(y') - \mathbb{1}(y' > y)]^2 dy', \quad (7)$$

286 where  $\mathbb{1}(y' > y)$  is the indicator function which equals zero when  $y' \leq y$  and equals one when  
 287  $y' > y$ . Another form of the CRPS which is more convenient to our purposes is the following :

$$\text{CRPS}(F, y) = \mathbb{E}_F |Y - y| - \frac{1}{2} \mathbb{E}_F |Y - Y'|, \quad (8)$$

288 where  $\mathbb{E}_F$  is the expectation and  $Y$  and  $Y'$  are random variables distributed according to  $F$ . The first  
 289 term of this equation is the mean-absolute error of the forecast. The second-term is the negative  
 290 half of the mean absolute difference between two variables distributed according to the forecast.  
 291 When we have an ensemble forecast, this is the mean absolute difference between two forecast  
 292 members. Using our notations, this allows to express easily the CRPS of the analog ensemble  
 293 forecast of  $y_i$  from the values  $\{y_j\}_{j \in I(i)}$  distributed according to  $\{p(j|i)\}_{j \in I(i)}$ .

$$\text{CRPS}_i = \sum_{j \in I(i)} p(j|i) y_{ji} - \frac{1}{2} \sum_{j, k \in I(i)} p(j|i) p(k|i) y_{jk}, \quad (9)$$

294 where we use the short notation for difference of absolute values  $y_{ji} := |y_j - y_i|$ . This can be  
 295 written more concisely as  $\text{CRPS}_i = \text{MAE}_i - \frac{1}{2} \text{MAD}_i$  with  $\text{MAE}_i := \sum_{j \in I(i)} p(j|i) y_{ji}$  and  $\text{MAD}_i :=$   
 296  $\sum_{j, k \in I(i)} p(j|i) p(k|i) y_{jk}$ . Finally, we are interested in optimizing:

$$\overline{\text{CRPS}} = \frac{1}{N} \sum_{i=1}^N \text{CRPS}_i. \quad (10)$$

297 Following similar steps as in the previous section, we find the following expression for the  
 298 gradient :

$$\frac{\partial \overline{\text{CRPS}}}{\partial \mathbf{A}} = \mathbf{A} \frac{1}{N} \sum_{i=1}^N \sum_{j \in I(i)} p(j|i) \left\{ \text{MAE}_i - \text{MAD}_i - \left( y_{ji} - \sum_{k \in I(i)} p(k|i) y_{kj} \right) \right\} \mathbf{x}_{ij} \mathbf{x}_{ij}^T. \quad (11)$$

299 This expression can be generalized to the case of vector-output  $\mathbf{y}$  by performing a sum of CRPS  
 300 over all coordinates.

301 If one is interested in computing the CRPS or of a multi-dimensional output, then one simple  
 302 possibility is to use the average of CRPS over all coordinates. The gradient of the coordinate-  
 303 averaged CRPS is then given by the above formula, replacing  $\text{MAE}_i$ ,  $\text{MAD}_i$  and  $y_{ij}$  terms by their  
 304 coordinate-average.

305 We provide modified formulas for the weighted-CRPS in appendix A, which could be used for  
 306 applications in which specific values (such as extreme values) are of interest.

#### 307 *d. Algorithm*

308 Goldberger et al. (2004) designed an algorithm which optimizes the matrix of a Mahalanobis-type  
 309 distance for classification purposes in a smoothed version of a nearest neighbour algorithm. Yang  
 310 et al. (2012) modified the algorithm of Goldberger et al. (2004), keeping only nearest neighbours  
 311 for classification. Weinberger and Tesauro (2007) designed an algorithm similar to Goldberger  
 312 et al. (2004), for optimizing a Mahalanobis-type distance for kernel regression. Here, we adapt  
 313 the algorithm from Weinberger and Tesauro (2007) to the case of analog forecasting, by keeping  
 314 only a finite number of nearest neighbours (analog) in the computation, and possibly adding  
 315 regularization terms. Therefore, our algorithm is a modified version of Weinberger and Tesauro

316 (2007) just as the algorithm of Yang et al. (2012) is a modified version of Goldberger et al. (2004).  
 317 Furthermore, we adapt the algorithm to the case where the loss function is not a mean-square error,  
 318 but a continuous ranked probability score, which is more suited to probabilistic ensemble-based  
 319 estimators used in atmospheric and ocean sciences.

320 One may be interested in sparse representation of the matrix  $\mathbf{A}$ , for instance to perform feature  
 321 selection as we demonstrate in section 4. To do so, we add the well-known  $l_1/l_2$  regularization  
 322 term to the loss function (Yin et al. 2014) :

$$\text{Loss}(\mathbf{A}) = \lambda \frac{\|\mathbf{A}\|_1}{\|\mathbf{A}\|_2} + \begin{cases} \text{MSE}(\mathbf{A}) \\ \text{CRPS}(\mathbf{A}) \end{cases} \quad (12)$$

323 with fixed parameter  $\lambda > 0$ , and  $\|\mathbf{A}\|_1$  and  $\|\mathbf{A}\|_2$  are the  $l_1$  and  $l_2$ -norms of matrix  $\mathbf{A}$ , namely the sum  
 324 of its coefficients' absolute values, and the square-root of the sum of its coefficients' squared values.  
 325 This additional term has sub-gradient  $\frac{\partial}{\partial \mathbf{A}} \frac{\|\mathbf{A}\|_1}{\|\mathbf{A}\|_2} = \frac{\text{sign}(\mathbf{A})}{\|\mathbf{A}\|_2} - \frac{\|\mathbf{A}\|_1}{\|\mathbf{A}\|_2^3} \mathbf{A}$ . This term has the advantage of  
 326 being scale-invariant, which is a desired property of our algorithm because adjusting the scale of  
 327  $\mathbf{A}$  allows to adjust the number of analogs used. If we had used a simple  $l_1$  Lasso regularization  
 328 (Tibshirani 1996) or  $l_2$  Ridge regularization (Hoerl and Kennard 1970), then we would have biased  
 329 our algorithm towards low-scale  $\mathbf{A}$ , *i.e.* towards a high number of analogs used. However, this  
 330 choice could be useful in applications where one wants to force the analog method to use a large  
 331 number of analogs, and therefore reducing variance at the cost of raising bias.

332 Finally, the matrix  $\mathbf{A}$  is updated at fixed learning rate  $\alpha > 0$  with the gradient descent rule (Bottou  
 333 2012):

$$\mathbf{A}_{new} = \mathbf{A}_{old} - \alpha \frac{\partial \text{Loss}}{\partial \mathbf{A}} \quad (13)$$

334 Based on that updated matrix  $\mathbf{A}_{new}$ , one can compute the new probabilities  $p(j|i)$  and associated  
 335 MSE and gradient, and update again the matrix using Eq. (13).

336 The algorithm is summarized in Algorithm 1. In this algorithm, the search for  $k$  neighbors is  
 337 denoted  $k\text{NN}$ . Particular cases of this algorithm include:

- 338 • *Weighting of coordinates*: by imposing  $\mathbf{A}$  to be diagonal.
- 339 • *Dimensionality reduction*: by setting  $\mathbf{A}$  to be of shape  $p \times d$  with  $p$  small.

340 Note that when one imposes  $\mathbf{A}$  to be diagonal, the tensor products  $\mathbf{x}_{ij}\mathbf{x}_{ij}^T$  must be replaced by  
 341 vectors of size  $d$  and whose coordinates are the square of the vector  $\mathbf{x}_{ij}$  (this is the diagonal of the  
 342 tensor  $\mathbf{x}_{ij}\mathbf{x}_{ij}^T$ ). This has the advantage of begin much less computationally-intensive than  $d \times d$   
 343 tensors.

---

**Algorithm 1:** Optimize Mahalanobis distance for Analog Prediction

---

**Function** OptimizeMahalanobis( $\mathbf{A}_0, \{\mathbf{x}_1, \dots, \mathbf{x}_N\}, \{y_1, \dots, y_N\}, \alpha, \lambda > 0, n, k$ )

```

   $\mathbf{A} = \mathbf{A}_0$ ;
  for  $t \in [1, n]$  do
    Build tree for  $\{\mathbf{A}\mathbf{x}_i, i \in [1, N]\}$ 
    for  $i \in [1, N]$  do
       $I(i) := k\text{NN}_j \|\mathbf{A}(\mathbf{x}_i - \mathbf{x}_j)\|$ ;
       $w_{ij} := \exp(-\|\mathbf{A}(\mathbf{x}_i - \mathbf{x}_j)\|^2)$ ;
       $p(j|i) := \frac{w_{ij}}{\sum_{k \in I(i)} w_{ik}}$ ;
    end
     $\mathbf{G} := (6) \text{ or } (11) \text{ or } (\text{A3})$ 
     $\mathbf{G} \leftarrow \mathbf{G} + \lambda \left( \frac{\text{sign}(\mathbf{A})}{\|\mathbf{A}\|_2} - \frac{\|\mathbf{A}\|_1}{\|\mathbf{A}\|_2^3} \mathbf{A} \right)$ 
     $\mathbf{A} \leftarrow \mathbf{A} - \alpha \mathbf{G}$ ;
     $E_t \leftarrow \begin{cases} \text{MSE}(\mathbf{A}) \\ \text{CRPS}(\mathbf{A}) \end{cases}$ 
  end
  return  $\mathbf{A}, (E_1, \dots, E_n)$ ;

```

---

344 **3. Lorenz System experiments**

345 To analyze the behaviour of our algorithm, we first use the well-known chaotic, three-variable  
 346 deterministic dynamical system of Lorenz (1963), in its usual setting, following the equations:

$$\frac{dx}{dt} = \sigma(y - x), \quad (14)$$

$$\frac{dy}{dt} = x(\rho - z) - y, \quad (15)$$

$$\frac{dz}{dt} = xy - \beta z, \quad (16)$$

347 with parameters  $\sigma = 10$ ,  $\beta = 8/3$  and  $\rho = 28$ . The equations are solved numerically using a  
 348 fourth-order Runge-Kutta explicit scheme (Butcher 1996) with time-step 0.01 (non-dimensional  
 349 time). This system of equation approximates the large-scale behaviour of atmospheric convection,



350 and bears properties of atmospheric and ocean circulation, in particular the sensitivity to initial  
351 condition.

352 Note that the classical notation  $(x, y, z)$  used here for the three coordinates of the system must  
353 not be confused with the notations  $\mathbf{x}$  and  $\mathbf{y}$  of the previous section, used to denote the predictor  
354 and predicted variables, respectively. In particular, this section will make the particular choice of  
355  $\mathbf{x} = [x(t), y(t), z(t)]$  (three coordinates as predictor) and either  $\mathbf{y} = [x(t+h), y(t+h), z(t+h)]$  (forecast  
356 of all three variables at horizon  $h$ ) or  $\mathbf{y} = [z(t+h)]$  (forecast of the last variable at horizon  $h$ ). For  
357 our numerical experiment, we choose to set  $\mathbf{A}$  to be a  $3 \times 3$  matrix, and therefore the transformed  
358 variables  $\mathbf{A}\mathbf{x}$  are of the same number (3) as the original variables  $\mathbf{x}$ . Note that other choices could  
359 have been retained for the shape of  $\mathbf{A}$ , such as a  $1 \times 3$  matrix to retain only one variable. We could  
360 also have included other input variables in the  $\mathbf{x}$  vector by using delayed-coordinates (Sauer et al.  
361 1991). The only constraint is that the number of columns in  $\mathbf{A}$  equals the size of  $\mathbf{x}$  (*i.e.* the number  
362 of input variables). We chose this simple  $3 \times 3$  setting with the original coordinates of the Lorenz  
363 system as input variables for illustrative testing of our method.

364 In all the experiments involving the Lorenz system, we will use  $k = 200$  analogs for each forecast  
365 ensemble, that is  $|I(\mathbf{x})| = 200$  for all  $\mathbf{x}$  (using notations from section 2). The set of analogs will  
366 be defined by the 200 indices  $j$  of nearest neighbors of  $\mathbf{x}$  according to the Euclidean distance  
367  $\|\mathbf{A}(\mathbf{x}_j - \mathbf{x})\|$ . We will use no regularization term in the algorithm (*i.e.*  $\lambda = 0$  in Eq. (12)), because  
368 it is unnecessary for this low-dimensional system.

369 The catalogs used in these experiments will be generated from long trajectories of the numerically-  
370 integrated Lorenz System. Elements of the catalog will comprise values of  $\mathbf{x}_i = [x(t_i), y(t_i), z(t_i)]$   
371 along with corresponding values of either  $\mathbf{y}_i = [x(t_i+h), y(t_i+h), z(t_i+h)]$  or  $\mathbf{y}_i = [z(t_i+h)]$  at  
372 horizon  $h > 0$ . The sequence of values  $\{t_1 \dots, t_N\}$  will be of the type  $t_i = t_1 + (i-1)\Delta t$  where  
373 the time separating two elements,  $\Delta t = 0.64$  (non-dimensional time), is chosen so that they can be  
374 considered independent. At least, this value of  $\Delta t = 0.64$  is enough so that  $x(t_{i+1})$  and  $x(t_{i-1})$  are  
375 not part of the 200 analogs of  $x(t_i)$  in our experiments.

376 Before applying our algorithm, we standardize each variable, that is we divide the values of  $x$  by  
377  $\langle (x - \langle x \rangle)^2 \rangle^{1/2}$ , the values of  $y$  by  $\langle (y - \langle y \rangle)^2 \rangle^{1/2}$ , and the values of  $z$  by  $\langle (z - \langle z \rangle)^2 \rangle^{1/2}$ . Therefore  
378 the matrix  $\mathbf{A}$  applies a transformation on standardized variables. For readability purposes, we will  
379 keep the notations  $x$ ,  $y$  and  $z$  for the standardized variables.

380 We will first apply our algorithm to minimize the RMSE of the analog ensemble average  
 381 prediction, and explore varying forecast horizons and forecasted variable, then we will investigate  
 382 varying catalog sizes, and finally we will compare the results of minimizing the RMSE versus  
 383 minimizing the CRPS of the analog forecast ensemble.

384 *a. Varying forecast horizon and variable*

385 In this subsection, we use a catalog of fixed size  $10^5$  generated from one long trajectory of the  
 386 Lorenz System, and we vary the forecast horizon  $h$  while keeping constant the values  $\{t_1, \dots, t_N\}$   
 387 as defined earlier.

388 To begin with, we apply our algorithm to minimize the MSE of the analog forecast at horizon  
 389  $h = 0.01$  (non-dimensional time), which is the time-step at which we perform the Runge-Kutta  
 390 integration of the Lorenz system’s equations. Note that at this very small horizon, the ideal forecast  
 391 is very close to a persistence forecast.

392 Our algorithm requires an initial value  $\mathbf{A}_0$  for the matrix  $\mathbf{A}$ . Note that as we are solving a  
 393 non-convex optimization problem, the choice of initial value influences the final result. We will  
 394 start with a standard choice for  $\mathbf{A}_0$  that we name “isotropic” ( $\mathbf{A}_0 = \mathbf{A}_{iso}$ ), that is the identity matrix:

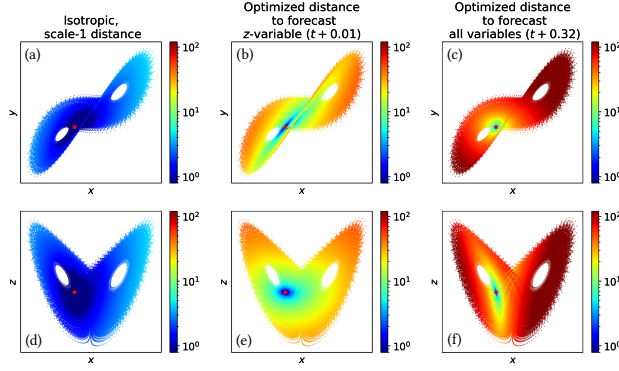
$$\mathbf{A}_{iso} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}. \quad (17)$$

395 Therefore, our algorithm starts with a baseline of standardized variables. Then, to choose the  
 396 learning rate  $\alpha$  (Eq. 13), we first compute the initial mean-squared error  $\text{MSE}_0$  of the analog  
 397 ensemble average, using distances  $\|\mathbf{A}_0(\mathbf{x}_i - \mathbf{x}_j)\|$  where  $\mathbf{A}_0$  is set to  $\mathbf{A}_{iso}$ , and the leave-one-out  
 398 methodology to compute the error, *i.e.* we make a forecast of  $\mathbf{y}_i$  using all the catalog without the  
 399 element  $(\mathbf{x}_i, \mathbf{y}_i)$ . Finally, we set the learning rate to be  $\frac{60}{\text{MSE}_0}$ , and we run our algorithm through  
 400  $n = 60$  iterations, after which the algorithm is converged to a value that we note  $\mathbf{A}_{con}(0.01)$  where  
 401 “con” stands for “converged”. Note that  $\text{MSE}_0$  depends on whether we are learning a matrix for  
 402 the forecast of the  $z$ -variable or of all three variables  $(x, y, z)$ . Fig. 2(a,b,d,e) allows to compare the  
 403 distances  $\|\mathbf{A}_{iso}(\mathbf{x}_i - \mathbf{x}_j)\|$  (Fig. 2.a,d) to the distances obtained after 60 iterations of our algorithm  
 404 for the forecast of the  $z$ -variable at horizon  $h = 0.01$  (Fig. 2.b,e). One can see that the optimized

405 distances are larger than the initial ones, this means that the probabilities given to the analogs  $p(i|j)$   
 406 are sharper: the selection of analogs is narrower. Also, one can witness a change in the relative  
 407 weights given to each coordinate: the ellipsis around the red dot has rotated, so that distances grow  
 408 faster with coordinate  $z$  in Fig. 2(e) than in Fig. 2(d) compared to how they grow with coordinate  
 409  $x$ . This shows that the optimal distance is anisotropic.

410 Then, the algorithm is run again, but at horizon  $h = 0.02$ , and starting from the previously  
 411 converged value  $\mathbf{A}_0 = \mathbf{A}_{con}(0.01)$ . Again, we compute the initial MSE at horizon  $h = 0.02$ , and  
 412 this time we set  $\alpha = \frac{30}{\text{MSE}_0}$  and run only  $n = 20$  iterations which is enough for the algorithm to  
 413 converge. This allows the algorithm to move to the closest local minimum of MSE when changing  
 414 slightly the forecast horizon. We then repeat this operation, raising the horizon by 0.01, running  
 415 the optimization with  $n = 20$  iterations after updating the value of  $\mathbf{A}_0$  and recomputing  $\alpha = \frac{30}{\text{MSE}_0}$   
 416 at each new horizon. Note that one succession of optimizations is run for the forecast of the  
 417 variable  $z$  and the other succession of optimizations is run independently for the forecast of the  
 418 three variables  $(x, y, z)$ . The result of running this procedure until horizon  $h = 0.32$  is shown in  
 419 Fig. 2(c,f). Comparing with Fig. 2(a,b,d,e) again shows a change in the orientation of the ellipsis  
 420 of constant distances, so that in particular the relative weights of each original coordinate in the  
 421 modified distance have changed. Also, all distances are larger indicating an even more selective  
 422 choice of analogs (*i.e.* even sharper distributions  $p(i|j)$ ).

423 Fig. 3(a,b) shows the evolution of all coefficients of the optimized  $\mathbf{A}_{conv}$  with forecast horizon,  
 424 and depending on the choice of either  $\mathbf{y}_i = [z(t_i + h)]$  or  $\mathbf{y}_i = [x(t_i + h), y(t_i + h), z(t_i + h)]$ . To aid the  
 425 interpretation, the  $l_2$ -norms of  $\mathbf{A}$ 's rows are also shown (*i.e.*  $\sqrt{A_{xj}^2 + A_{yj}^2 + A_{zj}^2}$  for  $j \in \{x, y, z\}$ ),  
 426 which allow to assess the relative importance of each original coordinate in the converged optimal  
 427 distance. This shows in particular that for small forecast horizons ( $h < 0.05$ ) the  $z$ -variable is the  
 428 most important variable when one is concerned with a forecast of the  $z$ -variable itself, which was  
 429 expected. Also, one can notice the growth of coefficients (and thus of the overall norm of  $\mathbf{A}$ ) with  
 430 horizon: this shows that a more selective choice of analogs must be done when concerned with a  
 431 forecast at large horizon, which is also expected due to the chaotic nature of the Lorenz system.  
 432 Note that the growth of the norm  $\mathbf{A}$ 's rows with horizon concerns the  $x$  and  $y$  coordinates. We  
 433 attribute this observation to the fact that these two coordinates are indicative of which “wing” one  
 434 is in (see Fig. 2): each “wing” corresponds to one of the two unstable fixed points around which

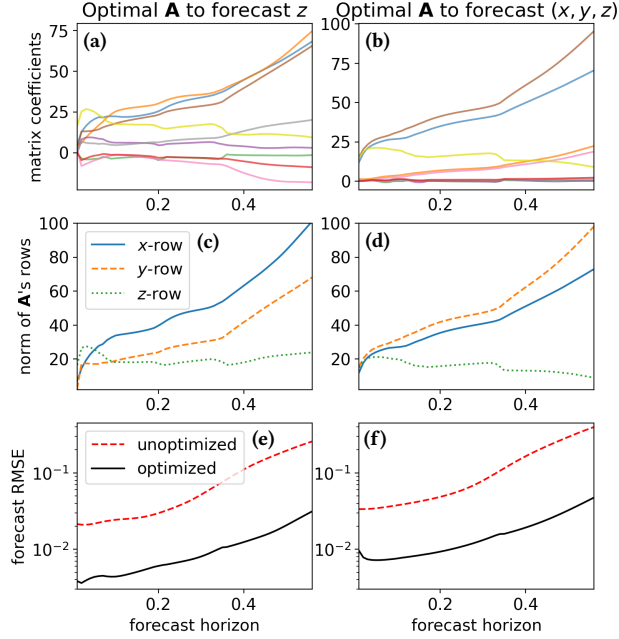


442 FIG. 2. Visualization of distances to a random point (red) across the Lorenz attractor, comparing (a,d)  
 443 unoptimized, isotropic distance with (b,e) distances optimized to minimize the RMSE of the mean analog  
 444 forecast of  $z$ -variable at short horizon and (c,f) mean analog forecast of all three state variables at large horizon.  
 445 The chosen forecast horizon can be interpreted as  $\sim 1$  hour and  $\sim 1$  day in atmospheric time scale.

435 solutions of the Lorenz system orbit before an eventual transition to the other wing. Finally, one  
 436 can see that even at large forecast horizon, the optimal distance found when trying to minimize the  
 437 RMSE of the forecast of only one variable differs from the optimal distance that is found when  
 438 trying to forecast all three variables of the Lorenz System. This indicates that there is no universal  
 439 distance that could be used for any analog method. Fig. 3(e,f) shows that the gain in RMSE after  
 440 optimization when compared to the RMSE found with  $\mathbf{A}_{iso}$  is constant throughout the range of  
 441 forecast horizons.

### 449 *b. Varying catalog size*

450 In this section, we will build catalogs of different sizes using the same long trajectory of  $6 \times 10^6$   
 451 non-dimensional time steps. From this trajectory, we will extract catalogs of different sizes,  
 452 and for each catalog size we extract 10 different catalogs to get variability in the results when  
 453 fixing only the catalog size. To do so, we first take regular sampling times  $\{t_1^*, \dots, t_{10N}^*\}$  of the  
 454 whole trajectory, separated by  $t_{i+1}^* - t_i^* = \frac{t_{10N}^* - t_1^*}{10N\Delta t}$  where  $N$  is the desired catalog size,  $\Delta t = 0.64$  (non-  
 455 dimensional time), and  $t_1^*$  and  $t_{10N}^*$  are the first and last time of the whole trajectory. Then, we take a  
 456 permutation  $\text{Perm}: [1, 10N] \mapsto [1, 10N]$ , and build 10 catalogs from indices  $\{t_{\text{Perm}(1)}, \dots, t_{\text{Perm}(N)}\}$ ,  
 457  $\{t_{\text{Perm}(N+1)}, \dots, t_{\text{Perm}(2N)}\}$ ,  $\{t_{\text{Perm}(2N+1)}, \dots, t_{\text{Perm}(3N)}\}$ , up to  $\{t_{\text{Perm}(9N+1)}, \dots, t_{\text{Perm}(10N)}\}$ . This  
 458 procedure allows to generate 10 catalogs from a single long trajectory for each chosen catalog size



446 FIG. 3. (a,b) Coefficients of matrix  $\mathbf{A}$  (see Fig. 4 for the color-to-coefficient correspondence) versus forecast  
 447 horizon. (c,d) Norm of rows of transform matrix  $\mathbf{A}$  versus forecast horizon. (e,f) Analog forecast RMSE versus  
 448 forecast horizon. (a,c,e) Forecast of  $z$ -variable. (b,d,f) Forecast of whole state-space vector.

459  $N$ . We use this procedure for 30 values of the catalog size  $N$  in the range  $N \in [4 \times 10^3, 4 \times 10^5]$   
 460 with regular sampling in log-scale within this interval.

461 In this experiment, we run our algorithm for the objective of minimizing the RMSE of the forecast  
 462 of all three variables at horizon  $h = 0.32$  (non-dimensional time). The algorithm was run for each  
 463 catalog constructed as explained above, with the linear transformation always initialized with the  
 464 identity matrix as  $\mathbf{A}_0 = \mathbf{A}_{iso}$ . The algorithm is run for 60 iterations, with a learning rate set to  
 465  $\frac{50}{\text{MSE}_0}$  (which depends on the catalog used) similarly to the previous experiment. The results of this  
 466 experiment are shown in Fig. 4. The ratio between optimized RMSE (with the final value of  $\mathbf{A}$ ) and  
 467 unoptimized RMSE (with  $\mathbf{A}_{iso}$ ) is nearly constant throughout the tested values of catalog size. The  
 468 values of coefficients of the final optimized matrix  $\mathbf{A}$  show a very strong dependency with catalog  
 469 size for small catalog sizes (below  $10^4$ ), also with greater variability of the coefficients for a fixed  
 470 catalog size. For larger values of the catalog size, the dependency of coefficients with catalog size  
 471 is much weaker, as well as the variability within a given value of catalog size. In particular, for  
 472 catalog sizes below  $\sim 10^4$ , the diagonal coefficients acting on the the  $x$  (blue lines) and  $z$  (yellow

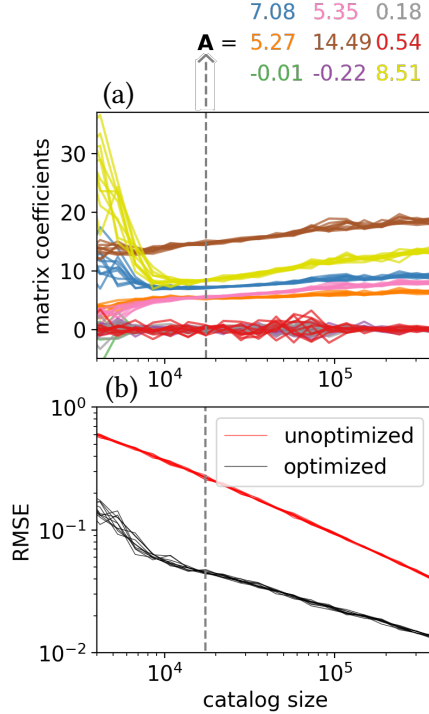
473 lines) variable are decreasing functions of catalog size, suggesting the need for a more rigorous  
474 selection of analogs when the catalog size is small. We interpret this result as the consequence  
475 of the fact that when the catalog size is small, the chance to find poor quality analogs in the 200  
476 nearest neighbors is higher, and therefore analog selection must be more meticulous. Moreover,  
477 this behaviour is not observed for the diagonal coefficient acting on  $y$ , which suggests that when  
478 changing the catalog size, one must not only adapt the number of analogs used, but also the very  
479 definition of the distance, such as the importance given to specific variables (or features) in the  
480 definition of the distance.

481 A striking property of Fig. 4(b) is the change of slope of the almost convex black curve of  
482 optimized RMSE around catalog size  $\sim 10^4$ . In particular, for catalog sizes above  $10^4$ , the black  
483 curve has a flatter slope than the red curve. This convexity of the black curve indicates that our  
484 optimization allows to approach faster the properties of an analog ensemble with large catalog size.  
485 In other words, the necessity to have a large catalog is less critical when optimizing the properties  
486 of analog methods. Our algorithm thus allows to (moderately) compensate for small-sized datasets.

487 Note that the behaviour observed for low values of the catalog size shows signs of the beginning  
488 of overfitting: coefficients are becoming large and highly variable. However, to check whether we  
489 are truly witnessing overfitting, we have computed the RMSE of the analog ensemble forecasts  
490 using the same catalogs as in the previous experiment, the same coefficients as fitted on these  
491 catalogs using our algorithm, but making forecasts on an independently generated test set of  $10^4$   
492 elements from a long trajectory, with elements separated by 0.64 non-dimensional time intervals.  
493 This independent test gives indistinguishingly (not shown) the same values as the RMSE shown in  
494 Fig. 4(b) and computed on the catalog (*i.e.* on the training set with leave-one-out methodology).  
495 This shows that the high values of the coefficients (1,1) and (3,3) of the fitted matrix  $\mathbf{A}$  really help  
496 to improve the performances of analog forecasting for small catalog size.

### 503 *c. Minimizing CRPS vs. minimizing RMSE*

504 Finally, we investigate the difference in the results of our algorithm when used with the same  
505 catalog, the same forecast objective (forecast of the  $z$  variable at horizon  $h = 0.04$ ), but with different  
506 loss function, either the RMSE or the CRPS.



497 FIG. 4. Matrix coefficients (a) of optimized linear transformation  $\mathbf{A}$  for average analog forecast of the Lorenz  
 498 system at horizon 0.32 (1 day), as a function of catalog size. RMSE is also shown (b) for both optimized (black)  
 499 and unoptimized (red) analog forecast. An example of optimized matrix  $\mathbf{A}$  is shown on top for catalog size 10  
 500 826. For each catalog size, 10 optimizations are run for 10 independent catalogs to account for variability in  
 501 the optimization process. Optimization is run at constant learning rate and fixed number of iterations, initialized  
 502 with the identity matrix. 200 analogs are retained for forecast.

507 We use the same catalog as in section 3.a), and we take the result of the experiment described  
 508 in the same section to obtain the matrix  $\mathbf{A}$  optimized for the RMSE of the average analog forecast  
 509 of the  $z$  variable at horizon  $h=0.04$ . To compare with the results for the CRPS, we use a similar  
 510 procedure to find the optimized matrix  $\mathbf{A}$ , starting by running our algorithm to minimize the CRPS  
 511 of the forecast of the  $z$ -variable at horizon 0.01, with  $\mathbf{A}_0 = \mathbf{A}_{iso}$ , 200 iterations, and a learning rate  
 512 of  $\frac{10^3}{\text{CRPS}_0}$ , where  $\text{CRPS}_0$  is defined analogously to  $\text{MSE}_0$  using the initial value of  $\mathbf{A}_0$ . Then, we  
 513 optimize the CRPS for the forecast of the  $z$ -variable at horizon  $h = 0.02$  initializing the matrix  
 514 with the previous optimized value  $\mathbf{A}_0 = \mathbf{A}_{con}(0.01)$  as was done for the RMSE. The algorithm is  
 515 run for 100 iterations at rate  $\frac{1.5 \times 10^3}{\text{CRPS}_0}$ . The operation is repeated to go from horizon 0.02 to 0.03, for  
 516 100 iterations at a rate  $\frac{1.5 \times 10^3}{\text{CRPS}_0}$ , and then for 40 iterations at a four times smaller rate to finalize the

517 convergence, a method called “step-decay schedule” for the learning rate (Ge et al. 2019). Finally,  
 518 to go up to horizon 0.04, we repeat the procedure with 60 iterations at rate  $\frac{1.5 \times 10^3}{\text{CRPS}_0}$  and then 40  
 519 iterations at a rate four times smaller.

520 One thing that we expect from our algorithm is that the analog ensemble with the distance  
 521 optimized to minimize the average CRPS would have a better representation of uncertainties.  
 522 Indeed, when optimizing the RMSE, it is not necessary that the ensemble spread corresponds  
 523 to the actual uncertainty of the forecast. In particular, other studies have reported RMSE-based  
 524 optimization of machine-learning algorithms to give over-confident ensemble forecasts (Frion et al.  
 525 2024). In our case, we observe the opposite behaviour: the RMSE-based optimization is under-  
 526 confident: the ensemble spread is larger than the actual uncertainty of the forecast. We interpret  
 527 this fact as the consequence of a linear properties of the Lorenz System at this small forecast  
 528 horizon and over the 200 selected analogs for this catalog size. In the case of linear dynamics, the  
 529 average value is closer to the truth when including more members.

530 To evaluate the goodness of uncertainty quantification between the two types of optimization, we  
 531 investigate first the prediction error, noted  $z_{ana}(t+h) - z_{truth}(t+h)$  in Fig. 5.a, where the subscript  
 532 “ana” refers to “analog”. This notation can be reconciled with the previous notations through the  
 533 following identity:

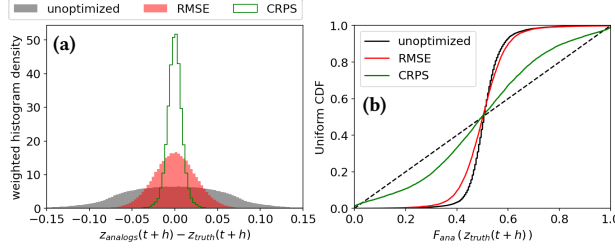
$$z_{ana}(t+h) - z_{truth}(t+h) = y_j - y_i. \quad (18)$$

534 This variable would be zero ideally. In particular, the probabilities  $p(j|i)$  should be highest  
 535 when this variable is close to zero. To evaluate this, we compute the empirical histograms of this  
 536 variable  $y_j - y_i$ , weighted by  $p(j|i)$ , for all values of  $i$  in the catalog and  $j \in I(i)$ , and comparing  
 537 the results of the three distances (unoptimized, optimized for RMSE, optimized for CRPS) in Fig.  
 538 5.a. We can see that all distributions are centered, and the CRPS is giving the sharpest distribution,  
 539 indicating a better uncertainty quantification.

540 Finally, in Fig. 5.b, we show the following variable:

$$F_{ana}(z_{truth}(t+h)) = \sum_{j \in I(i)} p(j|i) \mathbb{1}(y_j < y_i), \quad (19)$$





555 FIG. 5. Verifying statistics of analog forecasting ensembles for the case of the forecast of the  $z$ -variable of the  
 556 Lorenz system at horizon  $h=0.04$ , computed on  $10^4$  points. (a) Empirical distribution of the difference between the  
 557 analog forecast and the truth, for different distances : unoptimized (grey, full), optimized for RMSE minimization  
 558 (red, full) and optimized for CRPS minimization (green, empty). Each value of  $z_{ana}(t+h) - z_{truth}(t+h)$  is given  
 559 a weight  $p(\mathbf{x}_{ana}|\mathbf{x}_{truth})$  in the empirical density estimate. (b) Probability-probability plot (P-P plot) between the  
 560 reference cumulative probability distribution of a uniform random variable (vertical axis) and the cumulative  
 561 probability distribution of the analogue forecasts applied to the true value  $z_{truth}(t+h)$ . The dashed line indicates  
 562 what a perfect forecast distribution would give.

541 which is the empirical cumulative distribution of the analog forecast applied to the true outcome  
 542 value. Ideally, if the true outcome was drawn from the analog forecast distribution, the variable  
 543  $F_{ana}(z_{truth}(t+h))$  would be uniformly distributed, which would mean that the true outcome value  
 544 is exactly  $Q\%$  of the time above the  $Q$ -percentile of the analog ensemble forecast distribution. In  
 545 practice, the uncertainty estimation from the analog ensemble forecast distribution is not perfect,  
 546 however one aim of our CRPS-oriented algorithm is to optimize the distance in order to improve  
 547 uncertainty quantification from the analog ensemble. This is verified in Fig. 5.b which plots the  
 548 quantiles of  $F_{ana}(z_{truth}(t+h))$  versus the quantiles of a uniform distribution. In the ideal case,  
 549 the empirical line would lie along the diagonal. In this figure it is clear that the empirical results  
 550 from the RMSE-based optimization is closer to the diagonal line than the results with unoptimized  
 551 isotropic distance, and that the results are even more satisfying with the CRPS-based optimization.  
 552 However, further improvement is needed, as the figure still indicates an overdispersive ensemble:  
 553 the uncertainty estimate from the analog ensemble is greater than the true uncertainty of the  
 554 forecast.

## 563 4. Tropical cyclone intensity forecasting

### 564 a. Intensity forecasting

565 In this section we apply our algorithm to the study of the tropical cyclone, a meteorological  
566 system for which actual datasets, even at the global scale, are of a limited size. In such a context,  
567 we evaluate whether our method overcomes the typical dimensionality issues encountered in most  
568 tropical cyclone forecasting applications.

569 We use data from the International Best Track Archive for Climate Stewardship (IBTrACS, Knapp  
570 et al. 2010), a global compilation of best-track datasets from multiple international meteorological  
571 agencies which includes estimates of the storm location, intensity, and size in four geographical  
572 quadrants, on a six-hourly basis. IBTrACS suffers from spatio-temporal heterogeneities, especially  
573 concerning  $R_{max}$ , a critical parameter for the system dynamics (Avenas et al. 2024a), but it remains  
574 the most global dataset on tropical cyclones. In addition, statistical relationships have been recently  
575 developed to backup  $R_{max}$  estimates from better known parameters included in this global dataset  
576 (Avenas et al. 2023) and are used here.

577 For our experiment, we use only parts of the IBTrACS dataset. First, we restrict to the North-  
578 Atlantic ocean basin, to ensure a unified data treatment and definition from American agencies.  
579 Second, we use data from years 2003-2022, as we want to use a specific variable ( $R_{34}$ , see definition  
580 in appendix B) which is not included before 2003. Also, we consider solely tropical cyclones (TC)  
581 with a maximum wind speed  $V_{max}$  greater than 33m/s (see definition in appendix B) at least once in  
582 its lifetime. After this selection, we are left with 110 TCs re-sampled at 3h-time step (see appendix  
583 B).

584 The output variable to be estimated with analogs is defined as the difference in maximum  
585 sustained wind speed :

$$y := \Delta V_{max}(t, h) := V_{max}(t + h) - V_{max}(t) , \quad (20)$$

586 while the 15 input variables used to find and rate analogs are:

$$\begin{aligned}
\mathbf{x} := & \left\{ V_{max}(t), R_{max}^{IBT}(t), R_{34}(t), f_{Cor}(t), u_{trans}(t), \right. \\
& v_{trans}(t), R_{max}^{A23}(t), \frac{dV_{max}}{dt}(t), \frac{dR_{max}^{IBT}}{dt}(t), \\
& \frac{dR_{34}}{dt}(t), \frac{df_{Cor}}{dt}(t), \frac{du_{trans}}{dt}(t), \\
& \left. \frac{dv_{trans}}{dt}(t), \frac{dR_{max}^{A23}}{dt}(t), T_{18}(t) \right\}.
\end{aligned} \tag{21}$$

587 See appendix B for a description of these variables. As in the Lorenz-63 experiments, all  
588 variables are normalized by their standard-deviation. In the present work, we use our algorithm  
589 to minimize the CRPS of the analog ensemble forecast of  $y = \Delta V_{max}(t, h)$ , imposing the linear  
590 transform matrix  $\mathbf{A}$  to be square-diagonal. Note that we make use of regularization terms only in  
591 the next section. Our algorithm is used here for the purpose of weighting the (normalized) input  
592 variables in  $\mathbf{x}$ , with 15 weighting coefficients corresponding to each diagonal element of  $\mathbf{A}$ .

593 Note that we do not wish to compete with state-of-the-art forecasting algorithms as this would  
594 require a whole dedicated study, in particular for the definition of the dataset and input variables.  
595 Rather, our interest is in evaluating the behaviour of our algorithm on a reanalysis and reduced  
596 dataset describing a real-life physical problem. Also, we wish to demonstrate the ability of  
597 our algorithm to fine-tune the weighting of variables, without prior physical knowledge, for the  
598 challenging task of tropical cyclone intensity forecasting (Emanuel and Zhang 2016; Cangialosi  
599 et al. 2020).

600 From the 110 TCs in our dataset, we use 73 TCs for training (2/3 of the dataset) and 37 TCs for  
601 test. The train/test splitting is random, and we repeat the experiment with 10 random splittings to  
602 assess the sensitivity of the results to the splitting choice. In the training phase (*i.e.*, the optimisation  
603 of  $\mathbf{A}$ ), we use all the training dataset to evaluate the average CRPS and its gradient. We use a  
604 special type of leave-one-out methodology: to forecast  $\Delta V_{max}(t, h)$  of a given TC, we search for  
605 analogs in other TCs, but we also allow the use of analogs from the same TC only with a minimum  
606 separation of  $\pm 3$  days. This allows to raise the number of potential analogs and simplifies the  
607 algorithm structure. However, to assess the algorithm's performance on the test set, analogs are  
608 only searched within the training set, and therefore a given TC cannot be used as an analog of itself.

609 The good performances of the algorithm on the test set justify the reliability of this procedure (see  
 610 below).

611 First, we run our algorithm without regularization (*i.e.*  $\lambda = 0$ , see Eq. 12), and for forecast  
 612 horizons  $h = 12\text{h}, 24\text{h}, 36\text{h}, \dots, 120\text{h}$ . For each horizon  $h$ , the transformation matrix is initialized  
 613 with the identity matrix  $\mathbf{A} = \mathbf{A}_{iso}$ , and we use a constant learning rate of  $\frac{10}{\text{CRPS}_0}$  where  $\text{CRPS}_0$  is  
 614 defined as previously as the average CRPS on the training set with  $\mathbf{A} = \mathbf{A}_{iso}$ . The algorithm is run  
 615 for 50 iterations, which is enough to reach convergence. Note that the algorithm is run 10 times  
 616 for each horizon as we use 10 random train/test splittings. We note  $\text{CRPS}_{con}$  the average CRPS  
 617 obtained with  $\mathbf{A} = \mathbf{A}_{con}$  the converged matrix, and we define the CRPS gain as:

$$\% \text{CRPS}_{gain} := 100 \frac{\text{CRPS}_0 - \text{CRPS}_{con}}{\text{CRPS}_0}. \quad (22)$$

618 Note that this definition can be used for the average CRPS on the training set (using analogs  
 619 from the training set to forecast TCs in the training set) and for the average CRPS on the test set  
 620 (using analogs from the training set to forecast TCs in the test set). Ideally the two gains would  
 621 be nearly identical, which would indicate that the weights optimized on the training set generalize  
 622 well to the test set. This is confirmed in Fig. 6(b), where we show for each horizon  $h$  the median  
 623 and percentiles of the CRPS gain on the 10 random train/set splittings. The CRPS gain is slightly  
 624 higher on the training set, which is the sign of a slight overfitting, however the gains on the training  
 625 and test sets are similar for every horizon and follow the same tendency of a growth of the gain with  
 626 horizon. Gains are substantial, ranging from 7% to 20%, showing the interest of our methodology  
 627 compared to the use of a brute-force unoptimized distance.

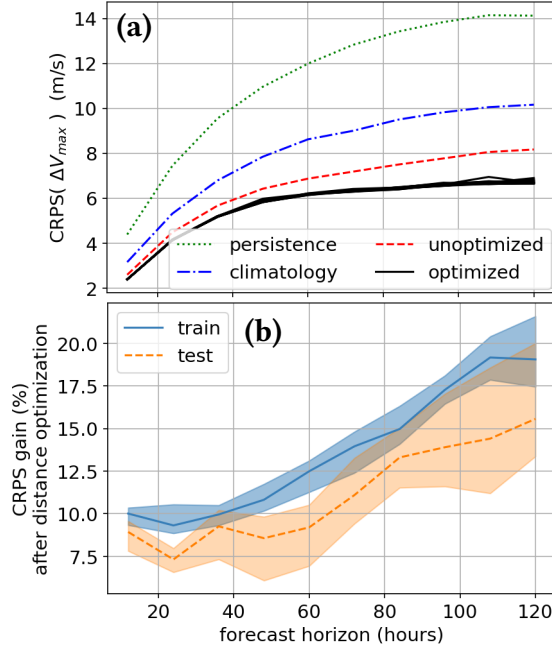
628 To assess the competitiveness of the analog methodology, we also compute the CRPS of two  
 629 benchmarks: persistence and climatological forecasts. Here, we define persistence forecast as a  
 630 deterministic forecast (*i.e.* a one-member ensemble) with  $\Delta V_{max} = 0$ , and its average CRPS is given  
 631 by its mean absolute error. We define the climatological forecast as an ensemble forecast where  
 632 all the elements of the dataset are used to build an ensemble forecast, and each element is given  
 633 equal weights. The climatological forecast is therefore given by the whole empirical distribution  
 634 of  $\Delta V_{max}$ , and therefore depends on the forecast horizon  $h$ . These benchmarks are evaluated on  
 635 the whole dataset (110 TCs), and compared to the (optimized and unoptimized) analog ensemble  
 636 forecasts on the whole dataset (train and test) using the leave-one-out methodology described

637 earlier. The corresponding average CRPS are shown as a function of horizon in Fig. 6(a). Note  
 638 that several black lines correspond to 10 different values of  $\mathbf{A}_{con}$  for each horizon  $h$ , associated  
 639 with the 10 random train/test splittings. The analog forecasts outperform the persistence and  
 640 climatological forecasts, especially for large forecast horizons, which is also where the CRPS gain  
 641 due to optimization is largest. The gain in CRPS for an horizon of 5 days thanks to our optimization  
 642 is of  $\sim 1.5\text{m/s}$ , which is close to the difference between the climatological forecast and the analog  
 643 forecast with unoptimized distance ( $\sim 2\text{m/s}$ ). Note that the climatological forecast can be viewed  
 644 as an analog ensemble forecast, where the number of analogs equals the size of the catalog, and  
 645 all distances are equal (*i.e.*  $p(i|j)$  is flat, using notations from section 2a). This means that our  
 646 algorithm allows for a gain in CRPS from a “naïve” (unoptimized) distance which is comparable  
 647 to the gain obtained when passing from a “flat” distance (climatology) to a naïve (unoptimized)  
 648 distance.

### 655 *b. Variable selection*

656 Then, we assess the ability of our algorithm to perform variable selection, using the regularization  
 657 term introduced in Eq. (12). We do so for forecast horizon  $h = 1$  day, and for regularization  
 658 coefficients  $\lambda = 0, 0.001, 0.002, 0.003, \dots, 0.015$ . For each value of  $\lambda$ , we take 10 random train/test  
 659 splittings as previously. We run our algorithm on each training set with constant learning rate  
 660 equal to  $\frac{10}{\text{CRPS}_0}$  where  $\text{CRPS}_0$  is defined as previously as the average CRPS on the training set with  
 661  $\mathbf{A} = \mathbf{A}_{iso}$ . The algorithm is run for 100 iterations this time, as the addition of the regularization  
 662 term requires more iterations to converge.

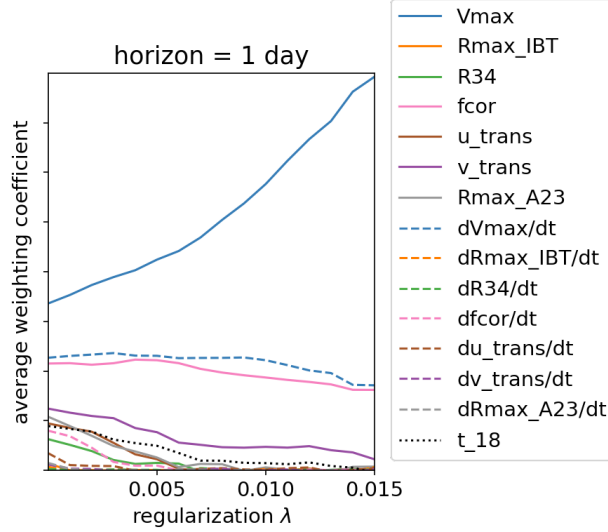
663 To see which variable is selected by the algorithm as we vary  $\lambda$ , we show in Fig. 7 the average  
 664 of each coefficient of  $\mathbf{A}_{con}$  over the 10 randomly selected training sets for each value of  $\lambda$ . This  
 665 shows that as  $\lambda$  grows, the smallest coefficients are drawn to zero, while the largest grow as a  
 666 compensation. For these experiments, the algorithm has selected  $V_{max}$ , as well as the Coriolis  
 667 frequency  $f_{Cor}$  although with a smaller weight. The meridional translation velocity is also selected  
 668 as a statistically relevant feature. The value of  $\frac{dV_{max}}{dt}$  is also statistically significant for the selection  
 669 of analogs, which seems reasonable as the local 3h-velocity growth rate is likely to be informative  
 670 of the short-term evolution.



649 FIG. 6. (a) Average CRPS (meters per second) of the forecast of  $\Delta V_{max}$  for different methods: persistence  
 650 (deterministic forecast with  $\Delta V_{max} = 0$ ), climatology (using the whole empirical distribution of  $\Delta V_{max}$  as an  
 651 ensemble), and unoptimized and optimized analog ensemble forecasts. (b) Gain in average CRPS of analog  
 652 forecast after optimizing the distance, on the training set used to optimize the distance (blue, full line) and on the  
 653 independent test set (orange, dashed line). The lines show the medians, while the shaded areas show the 25%  
 654 and 75% percentiles.

674 However Fig. 7 shows an average over 10 realizations of the training set, but the results for each  
 675 set can differ. Also, in practice, when doing variable selection, one is not interested in a particular  
 676 value of  $\lambda$  but in a fixed number of variables. To take this practical point of view, we take all  
 677 the values of converged  $\mathbf{A}_{con}$  from this experiment, and rank them by the number of coefficients  
 678 of  $\mathbf{A}$  which are above 0.15, an arbitrary threshold which gives an idea of the number of variables  
 679 selected by the algorithm. This allows to compute what is the average weight given to each variable  
 680 when fixing the number of non-negligible coefficients (Fig. 8(a)), as well as the probability to pick  
 681 a given variable when the number of non-negligible coefficients is fixed (Fig. 8(b)).

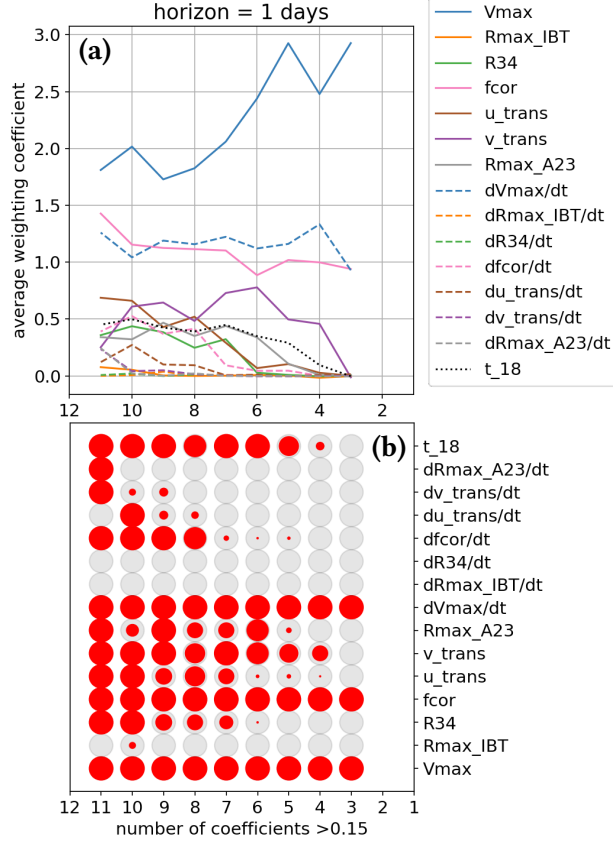
682 Although Fig. 8(a) mostly reproduces the behaviour of Fig. 7(a), some other features can be  
 683 extracted from this figure. One first striking fact is that the algorithm never selects more than  
 684 11 variables out of 15, even when the regularization parameter is set to zero. This shows that



671 FIG. 7. Weighting coefficients after optimization for the  $h = 1$  day-forecast of  $\Delta V_{max}$ , as a function of regular-  
 672 ization parameter  $\lambda > 0$ . The coefficients are averaged over 10 random splittings of the dataset into training and  
 673 test sets.

685 the algorithm is able to discard statistically irrelevant features without having to impose scarcity  
 686 through regularization. Also, we witness some nontrivial behaviour, in particular the fact that  
 687  $\frac{du_{trans}}{dt}$  is always selected when using 10 variables, but never when using 11. This goes against the  
 688 idea of the iterative algorithm of Alessandrini et al. (2018), described in the introduction. This  
 689 other algorithm goes from right to left in our Fig. 8, starting by adjusting the coefficients of the  
 690 most relevant variables used alone, and then adding new variables iteratively. Here, we show that  
 691 a variable which is relevant when used with a certain number of variables may not be relevant  
 692 when using a larger number of variables. From a methodological standpoint, this means that one  
 693 should in principle select all relevant variables at once and not iteratively, one-by-one as proposed  
 694 by Alessandrini et al. (2018).

695 Finally, note that the CRPS of the analog ensemble forecast with converged  $\mathbf{A}_{con}$  is a decreasing  
 696 function of  $\lambda$ : a higher regularization induces a larger error. For horizon  $h = 1$  day, the decrease  
 697 of  $CRPS(\mathbf{A}_{con})$  with  $\lambda$  is significant (at most a loss of 4.6% in  $CRPS_{gain}$  for  $\lambda = 0.015$ , while  
 698 the  $CRPS_{gain}$  has values  $\sim 9\%$ , not shown). We also do not witness a better generalization when  
 699 applying regularization: the ratio between the CRPS evaluated on the test set and the one evaluated  
 700 on the training set is not higher when raising  $\lambda$ . This shows that the number of variables selected by



703 FIG. 8. (a) Average weighing coefficient for each input variable, when aggregating converged values of  $\mathbf{A}_{con}$   
 704 for which the number of coefficients  $> 0.15$  is fixed. (b) The horizontal axis is the same as in the top panel, while  
 705 the vertical axis corresponds to the input variables. The radius of the red disks is proportional to the number of  
 706 times that this variable is selected with a coefficient  $> 0.15$  in all aggregated values of  $\mathbf{A}_{con}$ . When the red disk  
 707 is as large as the grey disk, this indicates that the variable is always selected. In contrary, when there is no red  
 708 disk in front of the grey disk, the variable is never selected.

701 the algorithm with  $\lambda = 0$  is already sufficiently small for the analog ensemble forecast to generalize  
 702 well, and imposing to use a lower number of variables is detrimental for forecast horizon  $h = 1$  day.

709 Note that several input variables are redundant here. For instance, the rate of change of the  
 710 Coriolis frequency  $\frac{df_{Cor}}{dt}$  can be expressed as a function of  $f_{Cor}$  and of the meridional translation  
 711 velocity. Also,  $R_{max}^{A23}$  is statistically determined from  $V_{max}$ ,  $f_{Cor}$  and  $R_{34}$ . The fact that these  
 712 variables are still selected some times by the algorithm indicates that creating new variables out  
 713 of existing ones for the definition of the distance may be useful in the case of analog forecasting,  
 714 in particular if these variables are non-linear functions of the initial variables that are based on



715 previous physical or empirical analysis showing their relevance. This is typically the case of  $R_{max}^{A23}$   
716 (see Eq. B3 in appendix B). Note also that this last variable is selected much more often than the  
717  $R_{max}^{IBT}$  from the IBTrACS dataset, although  $R_{max}^{A23}$  is a statistical approximation for the true radius  
718 of maximum wind speed. This confirms both the limitations of this parameter in the IBTrACS  
719 database (Combot et al. 2020) and the utility of empirical approaches to overcome this issue when  
720 studying the tropical cyclone dynamics (see also the discussion in Avenas et al. 2023). However,  
721 these experiments on IBTrACS data show that defining relevant variables is not enough, and should  
722 be complemented by an approach such as the one proposed here to systematically, at least, tune the  
723 weights given to each variable and the overall scale of the distances used to rate analogs.

## 724 **5. Conclusion and perspectives**

725 We have shown algorithms originally developed in the field of “distance learning”, which is a  
726 sub-field of machine-learning, can be adapted to allow for the optimization of the distance used in  
727 analog methods. To our knowledge, this is the first time that the gap between distance learning and  
728 analog methods is bridged. Our algorithm learns a linear transformation of the feature variables  
729 that is applied ahead of a classical Euclidean-distance-based analog ensemble methodology. This  
730 is equivalent to learning a “Mahalanobis-like” distance, but differs from using the Mahalanobis  
731 distance where the data’s covariance matrix is used directly. Distance learning algorithms were  
732 initially designed for classification purposes, with little interest in quantify uncertainties. On the  
733 contrary, our algorithm is designed for continuous estimation purposes (regression), such as the  
734 forecasting or downscaling of scalar variables. Furthermore, our algorithm tunes the distance so  
735 that the analog ensemble gives an accurate estimation of uncertainty while staying as close as  
736 possible to the ground truth. This is done in particular through the use of the continuous-ranked  
737 probability score as a loss function.

738 Our tests of the algorithm on analog forecasts of the three-variable chaotic Lorenz System show  
739 a non-trivial dependency of the optimal distance with forecast horizon, as well as catalog size. For  
740 low-size datasets, we observe strong variations of the optimal distance with catalog size, followed  
741 by stabilization above a given threshold and eventually convergence. We also observe the growth  
742 of the scale of the optimal distance with forecast horizon, indicating that a more severe selection  
743 of analogs is needed for long-term forecasts. These examples show that the optimal distance

744 strongly depends on the system under study, the objective task, and the number of available data.  
745 Finally, we show that a CRPS-based optimization allows to have better uncertainty quantification  
746 from analog ensembles compared to RMSE-based optimization that were developed previously in  
747 distance learning algorithms. This demonstrates the benefit of our adapted algorithm for the case  
748 of analog methods in atmospheric and ocean science.

749 To investigate the behaviour of our algorithm on a real system, we use IBTrACS tropical cyclone  
750 data, and test the ability of our algorithm to weight input variables in the case of intensity forecasting.  
751 First, analog methods outperform simple methods such as persistence or climatological forecasts at  
752 all forecast horizons in terms of CRPS. Second, our algorithm allows for significant improvement  
753 with respect to a baseline of analog forecast where all input variables are given equal weights.  
754 Third, even without regularization our algorithm already removes some irrelevant input variables,  
755 allowing for a first dimension reduction. To further reduce the number of variables used we add a  
756 regularization term, allowing to reveal which variables contribute the most to the optimal definition  
757 of distance for analog ensemble forecasting of tropical cyclone intensity. This demonstrates that  
758 our algorithm can be used on small-size datasets, which is an interesting property for TC studies,  
759 and that our algorithm allows to perform dimension reduction which is a key requirement of analog  
760 methods.

761 We note here that extensions of this algorithm were tried but not retained. These include the  
762 definition of a state-dependant distance, where the matrix  $\mathbf{A}$  is itself a (smooth) function of  $\mathbf{x}$ . This  
763 is one possible way of having a non-linear transformation, which is a generalization of our linear  
764 (constant) transformation  $\mathbf{A}$ . We have also tried to optimize several distance used at once, building  
765 several analog ensemble for each forecast, each with a different weight that is also optimized in the  
766 algorithm. Although feasible in practice, these extensions of our algorithm were computationally  
767 more intensive and did not yield remarkable improvements.

768 Our algorithm could still be modified in several ways to deal with existing problems in atmo-  
769 spheric and ocean science. First, note that the experiments conducted for this study were all  
770 performed on a personal laptop, but memory issues would arise in the case of both high number  
771 of features ( $>100$ ) and large number of training samples, which could be the consequence of using  
772 gridded fields of geophysical variables as features. These memory issues would be due to the prod-  
773 ucts  $\mathbf{x}_{ij} \mathbf{x}_{ij}^T$  which is as large as the square of the number of features (unless  $\mathbf{A}$  is diagonal). However,

774 there are memory-optimal ways to estimate such products, such as low-rank matrix approximations  
775 (Kumar et al. 2012). Also, in the experiments performed here we have used batch-gradient descent,  
776 which means that we use the whole training sample to compute the gradient (the whole sums over  
777  $N$  in sections 2.b and 2.c), but other techniques such as stochastic gradient descent (Bottou 2012) or  
778 mini-batch gradient descent (Khirirat et al. 2017), which compute the gradient over subsets of the  
779 training set, would allow to diminish the memory requirements of the method. These techniques  
780 could also help escaping sub-optimal local minima, since we are facing a non-convex optimization  
781 problem. More generally, routines available from machine-learning libraries allow to compute  
782 gradients very efficiently, which already enabled us to perform experiments on feature vectors of  
783 size exceeding  $10^3$  on a personal laptop (not shown).

784 The case of extreme events could be tackled by using our algorithm for weighted-CRPS mini-  
785 mization and giving more weights to the large values. This has potential applications for statistical  
786 downscaling of extreme precipitation, for instance.

787 Finally, note that there are numerous distance learning algorithms that are different from the  
788 ones we have used and could also be modified to meet the requirements of analog methods in  
789 ocean and atmospheric science. In particular, some algorithms have the property of solving convex  
790 optimization problems, including for instance Globerson and Roweis (2005). We are currently  
791 working on such adaptations.

792 *Acknowledgments.* P. Platzer, B. Chapron and A. Avenas acknowledge the support from  
 793 ERC project 856408-STUOD. A. Mouche acknowledges the support from ESA MAXSS  
 794 (4000132954/20/I-NB) and ESA MPC-S1 projects (4000107360/12/I-LG).

795 *Data availability statement.* The data used and generated for this article are available upon  
 796 request.

## 797 APPENDIX A

### 798 Generalization to weighted CRPS

799 An interesting feature of the CRPS is the possibility to weight the CRPS, and therefore give more  
 800 importance to specific outcome (*e.g.*, extreme values). The weighted CRPS is defined as:

$$\text{wCRPS}(F, y) = \int_{-\infty}^{\infty} [F(y') - \mathbb{1}(y' > y)]^2 w(y) dy', \quad (\text{A1})$$

801 for any non-negative function  $w(y)$ , and can be further expressed as (Taillardat et al. 2023):

$$\text{wCRPS}(F, y) = \mathbb{E}_F |W(Y) - W(y)| - \frac{1}{2} \mathbb{E}_F |W(Y) - W(Y')|, \quad (\text{A2})$$

802 where  $W(y) = \int_{-\infty}^y w(y') dy'$  is any primitive of  $w$ . Using our notations, rewriting  $W_{ij} :=$   
 803  $|\int_{y_i}^{y_j} w(y') dy'|$ , noting  $\text{wMAE}_i := \sum_{j \in I(i)} p(j|i) W_{ji}$  and  $\text{wMAD}_i := \sum_{j, k \in I(i)} p(j|i) p(k|i) W_{jk}$ , we  
 804 find that such a weighted CRPS has gradient:

$$\frac{\partial \overline{\text{wCRPS}}}{\partial \mathbf{A}} = \mathbf{A} \frac{1}{N} \sum_{i=1}^N \sum_{j \in I(i)} p(j|i) \left\{ \text{wMAE}_i - \text{wMAD}_i - \left( W_{ji} - \sum_{k \in I(i)} p(k|i) W_{kj} \right) \right\} \mathbf{x}_{ij} \mathbf{x}_{ij}^T. \quad (\text{A3})$$

## 805 APPENDIX B

### 806 Variables used for tropical cyclone forecasting

807 Variables used in this paper for the tropical cyclone forecasting experiment come directly or  
 808 indirectly from the IBTrACS database. IBTrACS is a compilation of the best-track data prepared  
 809 by the different Regional Specialized Meteorological Centers (RSMCs) and Tropical Cyclone  
 810 Warning Centers (TCWCs). Based on their area of responsibility, these regional agencies provide  
 811 analyses of the TC location, intensity and structure on a regular time basis using the available data.

812 IBTrACS variables used in this study are the following:

- 813 •  $V_{max}$ , the “maximum sustained wind speed”;
- 814 •  $R_{max}^{IBT}$  (we use the  $^{IBT}$ -superscript to distinguish from  $R_{max}^{A23}$  introduced below) is the radius  
815 of maximum sustained wind speed, defined as the distance between the TC center and the  
816 position at which  $V_{max}$  is measured;
- 817 •  $R_{34}$ , the radius at which the velocity reaches 34 knots (1 kt  $\approx$  0.51 m/s) in four geographical  
818 quadrants (NE, SE, SW, and NW);
- 819 •  $lat$ , the latitude of the TC center;
- 820 •  $storm\_speed$ , the storm translation speed;
- 821 •  $storm\_dir$ , the storm translation direction.

822 Because of varying definitions of the maximum sustained wind speed across the different agen-  
823 cies, we selected only USA agencies (*i.e* National Hurricane Center, Joint Typhoon Warning Center,  
824 and Central Pacific Hurricane Center) which all provide the 1-minute maximum sustained wind  
825 speed.

826 Furthermore, to focus on the strongest storms, and to ensure well-defined  $R_{34}$  values, we removed  
827 all storms with a lifetime maximum intensity lower than 17.5 m/s. We also considered storms for  
828 which  $R_{max}^{IBT}$  was defined for at least 72 consecutive hours. Lastly, we cropped all storm time series  
829 to select the part of each events for which  $V_{max}$  was comprised between 17.5 m/s and the lifetime  
830 maximum intensity, to investigate the intensification period.

831 In IBTrACS, some storm tracks are given on a six-hourly basis, while others are interpolated  
832 and thus given on a three-hourly basis. After applying the procedure mentioned above, the 111  
833 remaining storm tracks were all given on a three-hourly basis, except one, that was removed for  
834 consistency.

835 Then, the selected IBTrACS parameters have been directly used or transformed into other  
836 variables more relevant for the present study. The transformed variables include:

- 837 •  $R_{34}$ , whose nonzero values were averaged over the four geographical quadrants;
- 838 •  $f_{Cor}$ , the Coriolis frequency, defined as  $f_{Cor} = 2\Omega \sin(lat)$ , where  $\Omega = 7.292 \times 10^{-5} \text{ s}^{-1}$  is the  
839 Earth angular velocity and  $lat$  is the latitude of the TC center;

- 840 •  $u_{trans}$  and  $v_{trans}$ , the TC translation speed in the zonal and meridional directions, computed  
841 with *storm\_speed* and *storm\_dir*;
- 842 •  $T_{18}(t)$ , the number of hours after which  $V_{max}$  has reached 17.5 m/s. By definition, for each  
843 element of our dataset,  $T_{18}(t) > 0$ .
- 844 •  $R_{max}^{A23}$ , the radius of maximum wind speed estimated using the procedure described in Avenas  
845 et al. (2023) and detailed below.

846 The procedure to estimate  $R_{max}^{A23}$  can be summarized in three steps. An estimate of the TC  
847 maximum sustained wind speed that would correspond to an azimuthal average of the wind field,  
848 is first performed, using

$$V_{max,1D} = 0.6967V_{max} + 6.1992. \quad (B1)$$

849 Second, the absolute angular momentum that an air parcel loses between  $R_{34}$  and  $R_{max}$  is  
850 estimated with the statistical relationship

$$\frac{M_{max,1D}}{M_{34}} = 0.531 \exp\{-0.00214(V_{max,1D} - 17.5m/s) - 0.00314(V_{max,1D} - 17.5m/s)(\frac{1}{2}f_{Cor}R_{34})\}, \quad (B2)$$

851 where  $M_{34}$  is defined as  $M_{34} = R_{34} * 17.5m/s + \frac{1}{2}f_{Cor}R_{34}^2$ . Lastly,  $R_{max}^{A23}$  is estimated using the  
852 absolute angular momentum definition

$$R_{max}^{A23} = \frac{V_{max,1D}}{f_{Cor}} \left( \sqrt{1 + \frac{2fM_{max,1D}}{V_{max,1D}^2}} - 1 \right). \quad (B3)$$

## 853 **References**

- 854 Alessandrini, S., L. Delle Monache, C. M. Rozoff, and W. E. Lewis, 2018: Probabilistic prediction  
855 of tropical cyclone intensity with an analog ensemble. *Monthly Weather Review*, **146** (6), 1723–  
856 1744.
- 857 Alexander, R., and D. Giannakis, 2020: Operator-theoretic framework for forecasting nonlinear  
858 time series with kernel analog techniques. *Physica D: Nonlinear Phenomena*, **409**, 132–520.

- 859 Avenas, A., B. Chapron, A. Mouche, P. Platzer, and L. Vinour, 2024a: Revealing short-term  
860 dynamics of tropical cyclone wind speeds from satellite synthetic aperture radar. *Scientific*  
861 *Reports*, **14** (1), 12 808.
- 862 Avenas, A., A. Mouche, J. Knaff, X. Carton, and B. Chapron, 2024b: On the tropical cyclone  
863 integrated kinetic energy balance. *Geophysical Research Letters*, **51** (16), e2024GL108 327.
- 864 Avenas, A., A. Mouche, P. Tandeo, J.-F. Piolle, D. Chavas, R. Fablet, J. Knaff, and B. Chapron,  
865 2023: Reexamining the estimation of tropical cyclone radius of maximum wind from outer  
866 size with an extensive synthetic aperture radar dataset. *Monthly Weather Review*, **151** (12),  
867 3169–3189.
- 868 Bellet, A., A. Habrard, and M. Sebban, 2022: *Metric learning*. Springer Nature.
- 869 Benestad, R. E., 2010: Downscaling precipitation extremes: Correction of analog models through  
870 pdf predictions. *Theoretical and Applied Climatology*, **100**, 1–21.
- 871 Bessafi, M., A. Lasserre-Bigorri, C. Neumann, F. Pignolet-Tardan, D. Payet, and M. Lee-Ching-  
872 Ken, 2002: Statistical prediction of tropical cyclone motion: An analog–cliper approach.  
873 *Weather and forecasting*, **17** (4), 821–831.
- 874 Bonnardot, F., H. Quetelard, G. Jumaux, M.-D. Leroux, and M. Bessafi, 2019: Probabilistic  
875 forecasts of tropical cyclone tracks and intensities in the southwest indian ocean basin. *Quarterly*  
876 *Journal of the Royal Meteorological Society*, **145** (719), 675–686.
- 877 Bottou, L., 2012: Stochastic gradient descent tricks. *Neural Networks: Tricks of the Trade: Second*  
878 *Edition*, Springer, 421–436.
- 879 Butcher, J. C., 1996: A history of runge-kutta methods. *Applied numerical mathematics*, **20** (3),  
880 247–260.
- 881 Cangialosi, J. P., E. Blake, M. DeMaria, A. Penny, A. Latta, E. Rappaport, and V. Tallapragada,  
882 2020: Recent progress in tropical cyclone intensity forecasting at the national hurricane center.  
883 *Weather and Forecasting*, **35** (5), 1913–1922.

884 Chen, P., H. Yu, B. Brown, G. Chen, and R. Wan, 2016: A probabilistic climatology-based analogue  
885 intensity forecast scheme for tropical cyclones. *Quarterly Journal of the Royal Meteorological*  
886 *Society*, **142 (699)**, 2386–2397.

887 Combot, C., A. Mouche, J. Knaff, Y. Zhao, Y. Zhao, L. Vinour, Y. Quilfen, and B. Chapron,  
888 2020: Extensive high-resolution synthetic aperture radar (sar) data analysis of tropical cyclones:  
889 Comparisons with sfmr flights and best track. *Monthly Weather Review*, **148 (11)**, 4545 – 4563,  
890 <https://doi.org/10.1175/MWR-D-20-0005.1>.

891 Delle Monache, L., F. A. Eckel, D. L. Rife, B. Nagarajan, and K. Searight, 2013: Probabilistic  
892 weather prediction with an analog ensemble. *Monthly Weather Review*, **141 (10)**, 3498–3516.

893 Elliott, R., 1943: Studies of persistent regularities in weather phenomena. *Synoptic Weather Types*  
894 *of North America*.

895 Elsberry, R. L., and H.-C. Tsai, 2014: Situation-dependent intensity skill metric and intensity  
896 spread guidance for western north pacific tropical cyclones. *Asia-Pacific Journal of Atmospheric*  
897 *Sciences*, **50**, 297–306.

898 Emanuel, K., and F. Zhang, 2016: On the predictability and error sources of tropical cyclone  
899 intensity forecasts. *Journal of the Atmospheric Sciences*, **73 (9)**, 3739–3747.

900 Fetanat, G., A. Homaifar, and K. R. Knapp, 2013: Objective tropical cyclone intensity estimation  
901 using analogs of spatial features in satellite data. *Weather and forecasting*, **28 (6)**, 1446–1459.

902 Fraedrich, K., C. C. Raible, and F. Sielmann, 2003: Analog ensemble forecasts of tropical cyclone  
903 tracks in the australian region. *Weather and forecasting*, **18 (1)**, 3–11.

904 Fraedrich, K., and B. Rückert, 1998: Metric adaption for analog forecasting. *Physica A: Statistical*  
905 *Mechanics and its Applications*, **253 (1-4)**, 379–393.

906 Frion, A., L. Drumetz, G. Tochon, M. D. Mura, and A. A. E. Bey, 2024: Koopman ensembles for  
907 probabilistic time series forecasting. *arXiv preprint arXiv:2403.06757*.

908 Ge, R., S. M. Kakade, R. Kidambi, and P. Netrapalli, 2019: The step decay schedule: A near  
909 optimal, geometrically decaying learning rate procedure for least squares. *Advances in neural*  
910 *information processing systems*, **32**.



911 Globerson, A., and S. Roweis, 2005: Metric learning by collapsing classes. *Advances in neural*  
912 *information processing systems*, **18**.

913 Goldberger, J., G. E. Hinton, S. Roweis, and R. R. Salakhutdinov, 2004: Neighbourhood compo-  
914 nents analysis. *Advances in neural information processing systems*, **17**.

915 Hastie, T., 2009: The elements of statistical learning: data mining, inference, and prediction.  
916 Springer.

917 Held, M., P. Wolfe, and H. P. Crowder, 1974: Validation of subgradient optimization. *Mathematical*  
918 *programming*, **6**, 62–88.

919 Hersbach, H., 2000: Decomposition of the continuous ranked probability score for ensemble  
920 prediction systems. *Weather and Forecasting*, **15 (5)**, 559–570.

921 Hersbach, H., and Coauthors, 2020: The era5 global reanalysis. *Quarterly Journal of the Royal*  
922 *Meteorological Society*, **146 (730)**, 1999–2049.

923 Hoerl, A. E., and R. W. Kennard, 1970: Ridge regression: Biased estimation for nonorthogonal  
924 problems. *Technometrics*, **12 (1)**, 55–67.

925 Horton, P., M. Jaboyedoff, and C. Obled, 2017: Global optimization of an analog method by means  
926 of genetic algorithms. *Monthly Weather Review*, **145 (4)**, 1275–1294.

927 Jackson, C., T. Ruff, J. Knaff, A. Mouche, and C. Sampson, 2021: Chasing cyclones from space.  
928 *Eos*, **102**, –, <https://doi.org/10.1029/2021EO159148>.

929 Jézéquel, A., P. Yiou, and S. Radanovics, 2018: Role of circulation in european heatwaves using  
930 flow analogues. *Climate dynamics*, **50 (3-4)**, 1145–1159.

931 Khirirat, S., H. R. Feyzmahdavian, and M. Johansson, 2017: Mini-batch gradient descent: Faster  
932 convergence under data sparsity. *2017 IEEE 56th Annual Conference on Decision and Control*  
933 *(CDC)*, IEEE, 2880–2887.

934 Kingma, D. P., and J. Ba, 2014: Adam: A method for stochastic optimization. *arXiv preprint*  
935 *arXiv:1412.6980*.

936 Knapp, K. R., M. C. Kruk, D. H. Levinson, H. J. Diamond, and C. J. Neumann, 2010: The  
937 international best track archive for climate stewardship (ibtracs) unifying tropical cyclone data.  
938 *Bulletin of the American Meteorological Society*, **91 (3)**, 363–376.

939 Krick, I. P., 1942: *A Dynamical Theory of the Atmospheric Circulation and Its Use in Weather*  
940 *Forecasting...: Studies of Persistent Regularities in Weather Phenomena*. California Institute of  
941 Technology.

942 Kumar, S., M. Mohri, and A. Talwalkar, 2012: Sampling methods for the nyström method. *The*  
943 *Journal of Machine Learning Research*, **13 (1)**, 981–1006.

944 Langmack, H., K. Fraedrich, and F. Sielmann, 2012: Tropical cyclone track analog ensemble  
945 forecasting in the extended australian basin: Nwp combinations. *Quarterly Journal of the Royal*  
946 *Meteorological Society*, **138 (668)**, 1828–1838.

947 Le Bras, P., F. Sévellec, P. Tandeo, J. Ruiz, and P. Ailliot, 2024: Selecting and weighting dynamical  
948 models using data-driven approaches. *Nonlinear Processes in Geophysics*, **31 (3)**, 303–317.

949 Lewis, W. E., T. L. Olander, C. S. Velden, C. Rozoff, and S. Alessandrini, 2021: Analog ensemble  
950 methods for improving satellite-based intensity estimates of tropical cyclones. *Atmosphere*,  
951 **12 (7)**, 830.

952 Lguensat, R., P. Tandeo, P. Ailliot, M. Pulido, and R. Fablet, 2017: The analog data assimilation.  
953 *Monthly Weather Review*, **145 (10)**, 4093–4107.

954 Lorenz, E. N., 1956: *Empirical orthogonal functions and statistical weather prediction*, Vol. 1.  
955 Massachusetts Institute of Technology, Department of Meteorology Cambridge.

956 Lorenz, E. N., 1963: Deterministic nonperiodic flow. *Journal of atmospheric sciences*, **20 (2)**,  
957 130–141.

958 Lorenz, E. N., 1969: Atmospheric predictability as revealed by naturally occurring analogues.  
959 *Journal of Atmospheric Sciences*, **26 (4)**, 636–646.

960 Lucarini, V., and Coauthors, 2016: *Extremes and recurrence in dynamical systems*. John Wiley &  
961 Sons.

- 962 Matulla, C., X. Zhang, X. Wang, J. Wang, E. Zorita, S. Wagner, and H. Von Storch, 2008:  
963 Influence of similarity measures on the performance of the analog method for downscaling daily  
964 precipitation. *Climate Dynamics*, **30**, 133–144.
- 965 McDermott, P. L., and C. K. Wikle, 2016: A model-based approach for analog spatio-temporal  
966 dynamic forecasting. *Environmetrics*, **27** (2), 70–82.
- 967 McLachlan, G. J., 1999: Mahalanobis distance. *Resonance*, **4** (6), 20–26.
- 968 Neumann, C. J., and J. R. Hope, 1972: Performance analysis of the hurran tropical cyclone forecast  
969 system. *Monthly Weather Review*, **100** (4), 245–255.
- 970 Nicolis, C., 1998: Atmospheric analogs and recurrence time statistics: Toward a dynamical  
971 formulation. *Journal of the atmospheric sciences*, **55** (3), 465–475.
- 972 Peterson, L. E., 2009: K-nearest neighbor. *Scholarpedia*, **4** (2), 1883.
- 973 Platzer, P., P. Yiou, P. Naveau, J.-F. Filipot, M. Thiébaud, and P. Tandeo, 2021a: Probability  
974 distributions for analog-to-target distances. *Journal of the Atmospheric Sciences*, **78** (10), 3317–  
975 3335.
- 976 Platzer, P., P. Yiou, P. Naveau, P. Tandeo, J.-F. Filipot, P. Ailliot, and Y. Zhen, 2021b: Using local  
977 dynamics to explain analog forecasting of chaotic systems. *Journal of the Atmospheric Sciences*,  
978 **78** (7), 2117–2133.
- 979 Sauer, T., J. A. Yorke, and M. Casdagli, 1991: Embedology. *Journal of statistical Physics*, **65**,  
980 579–616.
- 981 Taillardat, M., A.-L. Fougères, P. Naveau, and R. De Fondeville, 2023: Evaluating probabilistic  
982 forecasts of extremes using continuous ranked probability score distributions. *International*  
983 *Journal of Forecasting*, **39** (3), 1448–1459.
- 984 Tandeo, P., and Coauthors, 2015: Combining analog method and ensemble data assimilation:  
985 application to the lorenz-63 chaotic system. *Machine Learning and Data Mining Approaches*  
986 *to Climate Science: proceedings of the 4th International Workshop on Climate Informatics*,  
987 Springer, 3–12.

- 988 Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of cmip5 and the experiment  
989 design. *Bulletin of the American meteorological Society*, **93** (4), 485–498.
- 990 Tibshirani, R., 1996: Regression shrinkage and selection via the lasso. *Journal of the Royal*  
991 *Statistical Society Series B: Statistical Methodology*, **58** (1), 267–288.
- 992 Toth, Z., 1991: Intercomparison of circulation similarity measures. *Monthly weather review*,  
993 **119** (1), 55–64.
- 994 Tsai, H.-C., and R. L. Elsberry, 2014: Applications of situation-dependent intensity and intensity  
995 spread predictions based on a weighted analog technique. *Asia-Pacific Journal of Atmospheric*  
996 *Sciences*, **50**, 507–518.
- 997 Tsai, H.-C., and R. L. Elsberry, 2019: Combined three-stage 7-day weighted analog intensity  
998 prediction technique for western north pacific tropical cyclones: Demonstration of optimum  
999 performance. *Weather and Forecasting*, **34** (6), 1979–1998.
- 1000 Van den Dool, H., 1994: Searching for analogues, how long must we wait? *Tellus A*, **46** (3),  
1001 314–324.
- 1002 Weickmann, L., 1924: Wellen im luftmeer. *Treatise in Math.-Phys. of the Saxon Academy of*  
1003 *Science*, **39** (2).
- 1004 Weinberger, K. Q., and G. Tesauro, 2007: Metric learning for kernel regression. *Artificial intelli-*  
1005 *gence and statistics*, PMLR, 612–619.
- 1006 Yang, W., K. Wang, and W. Zuo, 2012: Fast neighborhood component analysis. *Neurocomputing*,  
1007 **83**, 31–37.
- 1008 Yin, P., E. Esser, and J. Xin, 2014: Ratio and difference of  $l_1$  and  $l_2$  norms and sparse rep-  
1009 resentation with coherent dictionaries. *Communications in Information and Systems*, **14** (2),  
1010 87–109.
- 1011 Yiou, P., 2014: Anawege: a weather generator based on analogues of atmospheric circulation.  
1012 *Geoscientific Model Development*, **7** (2), 531–543.

- 1013 Yiou, P., M. Boichu, R. Vautard, M. Vrac, S. Jourdain, E. Garnier, F. Fluteau, and L. Menut, 2014:  
1014 Ensemble meteorological reconstruction using circulation analogues of 1781–1785. *Climate of*  
1015 *the Past*, **10** (2), 797–809.
- 1016 Yiou, P., and A. Jézéquel, 2020: Simulation of extreme heat waves with empirical importance  
1017 sampling. *Geoscientific Model Development*, **13** (2), 763–781.
- 1018 Zhao, Z., and D. Giannakis, 2016: Analog forecasting with dynamics-adapted kernels. *Nonlinear-*  
1019 *ity*, **29** (9), 2888.
- 1020 Zhen, Y., P. Tandeo, S. Leroux, S. Metref, T. Penduff, and J. Le Sommer, 2020: An adaptive  
1021 optimal interpolation based on analog forecasting: application to ssh in the gulf of mexico.  
1022 *Journal of Atmospheric and Oceanic Technology*, **37** (9), 1697–1711.
- 1023 Zorita, E., and H. Von Storch, 1999: The analog method as a simple statistical downscaling  
1024 technique: Comparison with more complicated methods. *Journal of climate*, **12** (8), 2474–2489.